



# Verdien av ekspertise

En empirisk studie på relasjonene ved menneskelig input i store språkmodeller ved verdivurdering av oppstartsselskaper, og interaksjonseffekten av ekspert- og ikke-ekspertinput.

**Casper M.B. Olsvik & Marcus J.V. Daae Lampe**

**Veileder: Eirik Sjøholm Knudsen**

Masterutredning økonomi og administrasjon

Hovedprofiler: finansiell økonomi & strategi og ledelse

NORGES HANDELSHØYSKOLE

Dette selvstendige arbeidet er gjennomført som ledd i masterstudiet i økonomi- og administrasjon ved Norges Handelshøyskole og godkjent som sådan. Godkjenningen innebærer ikke at Høyskolen eller sensorer inntår for de metoder som er anvendt, resultater som er fremkommet eller konklusjoner som er trukket i arbeidet.

## Forord og anerkjennelser

Denne oppgaven markerer kulminasjonen av vår mastergrad i økonomi og administrasjon ved Norges Handelshøyskole. Den er et resultat av en spennende og tidvis utfordrende faglig reise. Arbeidet med oppgaven har gitt oss en mulighet til å benytte mye tid på å fordype oss i noe vi finner interessant. Vi har med avhandlingen hatt som formål å bedre forstå hvordan store språkmodeller (LLM, KI) påvirkes av ulik input. Vi håper vår forskning vil bidra til at KI benyttes med større bevissthet rundt dens muligheter og begrensninger.

Forfatterne er Marcus J. V. Daae Lampe med hovedfag finansiell økonomi og Casper M. B. Olsvik med hovedfag strategi og ledelse. Oppgaven har ført oss gjennom utfordringer, mestring og lange dager. En viktig faktor for oss i denne prosessen har vært støtten fra de rundt oss, venner og familie, det setter vi pris på.

Vi vil uttrykke en spesiell takk til vår veileder Professor Eirik Sjøholm Knudsen. Eirik har bidratt med verktøy, tilbakemeldinger og uvurderlig innsikt helt fra formulering av forskningsspørsmål til innlevering. Vi vil også rette en takknemlighet til Norges Handelshøyskole som institusjon, med de ressurser, engasjerte forelesere og et trygt og godt læringsmiljø som har hjulpet oss å utvikle oss akademisk. Til slutt vil vi takke våre medstudenter som har gjort den akademiske reisen til noe mer enn bare en akademisk reise.

Takk.

## Sammendrag

Tidligere studier har vist at kvaliteten på menneskelig input kan forbedre interaksjon og feedback fra store språkmodeller (LLM). Dette indikerer at språkmodellene er mottagelige for menneskelig input. Med dette som utgangspunkt undersøkte denne studien om opplevd ekspertise i input påvirket LLMs verdivurdering av startups. I tillegg undersøkte studien om det var systematiske forskjeller mellom de ledende språkmodellene ChatGPT, Claude, Llama og Gemini. Ved å gjøre et utvalg av 50 norske startups, presenterte vi selskapenes prospekter ved kapitalinnhenting til språkmodellene og spurte om en investeringsanbefaling og en verdivurdering. Språkmodellene ble deretter gitt input som gikk imot deres opprinnelige vurdering. Denne inputen ble presentert som om den kom enten fra en ekspert eller en ikke-ekspert på feltet, med fokus på verdivurdering av oppstartsselskaper. Resultatene indikerte at alle modellene endret sine opprinnelige vurderinger basert på menneskelig input. Selv om Gemini var mer skeptisk mot å endre mening uten ytterligere opplysning. Det var ingen signifikant forskjell på hvordan modellene ble påvirket av ekspert- eller ikke-ekspertinput når modellene ble undersøkt samlet. Dette viser at store språkmodeller kan generelt være sensitiv til uenighet, uavhengig av kilden. Studien setter søkelys på behovet for ytterligere forskning på hvordan forbedre robuste og pålitelige språkmodeller ved integrasjon av menneskelig input. Spesielt som beslutningsstøtte i situasjoner med høy usikkerhet og risiko som i Venture Capital (VC)-bransjen.

**Nøkkelord:** Store språkmodeller (LLMs), verdivurdering, oppstartsselskaper, ekspertise, menneskelig input, venture capital, bias.

---

# Innholdsfortegnelse

<b>Forord og anerkjennelser</b> .....	2
<b>Sammendrag</b> .....	3
<b>1. Introduksjon</b> .....	6
1.1 Bakgrunn for studien.....	6
1.2 Forskningsspørsmål .....	7
1.3 ChatGPT forventinger og faktiske innsikter fra studien.....	8
1.4 Formål og struktur.....	9
<b>2. Teoretisk rammeverk</b> .....	11
2.1 Store språkmodeller (LLMs): En introduksjon .....	11
2.1.1 Hva er LLMs og hvordan fungerer de?.....	11
2.1.2 Bruken av LLMs i verdsettelse av selskaper .....	12
2.2 Startup-verdivurdering og Venture Capital-modeller.....	13
2.2.1 Tradisjonelle verdivurderingsmetoder.....	13
2.2.2 Verdivurdering av startups.....	15
2.3 Kunstig Intelligens i Venture Capital-bransjen .....	16
2.3.1 Potensialet for KI i beslutningsstøtte .....	16
2.3.2 Tidligere forskning rundt KI og VC-beslutninger .....	17
2.4 Menneskelig input og store språkmodeller .....	19
2.4.1 Hvordan språkmodeller justerer seg basert på menneskelig input .....	19
2.4.2 Forskjeller mellom ekspert- og ikke-ekspertinput .....	19
2.5 Hypotese-formulering .....	20
<b>3. Metode</b> .....	22
3.1 Forskningsmodell.....	22
3.2 Forskningsdesign.....	22
3.3 Datainnsamling .....	23
3.3.1 Utvalgsprosedyre.....	23
3.3.2 Undersøkelsen .....	24
3.3.3 Forberedelse av data.....	28
3.3.4 Utvalgets egenskaper .....	30
3.4 Målinger .....	31
3.4.1 Avhengig variabel.....	31
3.4.2 Uavhengige variabler.....	32
3.4.3 Kontrollvariabler .....	32

---

3.5 Dataanalyse .....	33
3.5.1 Beskrivende statistikk .....	33
3.5.2 Kjikvadrattest .....	34
3.5.3 Logistisk regresjonsanalyse .....	34
3.5.4 T-test .....	35
3.5.5 Lineær regresjonsanalyse .....	36
3.5.6 Analyse av påvirkningsgrad .....	36
3.5.7 Analyse av ulike LLMs .....	37
3.6 Evaluering av validitet og reliabilitet .....	37
3.6.1 Indre validitet .....	38
3.6.2 Ytre validitet .....	39
3.6.3 Reliabilitet .....	40
3.7 Etikk .....	41
<b>4. Resultater</b> .....	<b>42</b>
4.1 Beskrivende statistikk .....	42
4.2 Test av hypotese 2 og hypotese 3 .....	48
4.3 Test av hypotese 1 .....	50
4.4 Test av hypotese 4 .....	53
4.5 Test av hypotese 4 med kontrollvariabler .....	55
4.5.1 Test av hypotese 4 med kontrollvariabler gitt positiv retning .....	61
4.5.2 Test av hypotese 4 med kontrollvariabler gitt negativ retning .....	66
<b>5. Diskusjon</b> .....	<b>72</b>
5.1 Hovedfunn .....	72
5.2 Ytterligere funn .....	75
5.3 Teoretiske implikasjoner .....	77
5.4 Praktiske implikasjoner .....	78
5.5 Forfatterne om fremtidige studier .....	78
5.6 "Selvinnsikt" ChatGPT 4o .....	79
5.7 Begrensinger ved studien .....	81
<b>6. Konklusjon</b> .....	<b>83</b>
<b>7. KI Erklæring</b> .....	<b>85</b>
<b>8. Referanser</b> .....	<b>86</b>

# 1. Introduksjon

## 1.1 Bakgrunn for studien

Studiens mål er å øke forståelsen for kunstig intelligente språkmodeller. De fleste store språkmodeller er i dag en såkalt black box. Som her vil si at vi ikke vet hvordan de fungerer, tar avgjørelser og reagerer på input fra mennesker (University of Michigan-Dearborn, 2023).

Om du spør ChatGPT hva du kan spise til middag så vet vi ikke konkret hvordan den kommer frem til oppskriften den sender til deg. Vi vet at den trenes på svært store datamengder, og at den kan benytte internett for å svare på spørsmål, men ikke hvordan.

KI kommer med stor sannsynlighet til å være med oss i økende grad i årene fremover. En studie gjennomført av Felten et al. (2023) konkluderte med at KI overlapper, altså kan utføre noen av oppgavene i større eller mindre grad til 980 av 1016 jobbkategorier undersøkt. Det kan være en risiko for samfunnet å lene seg på systemer vi ikke vet helt hvordan fungerer. Spesielt hvilke bias som finnes i treningsdataene til språkmodellene. Språkmodellene er trent på menneskelig data, og mennesker er utsatt for kognitiv bias. Kahneman et al. (1982) definerer dette som underbevisste og systematiske feil i tenkning som skjer når mennesker prosesserer og tolker informasjon i deres omgivelser og påvirker deres beslutninger og vurderinger. Språkmodellene trenes på store datasett produsert av mennesker. Dette får konsekvenser for hvordan språkmodellene reagerer på input fra oss mennesker i form av prompts i samtaler. Annen litteratur har funnet at språkmodellene justerer sine evalueringer basert på menneskelig input. Språkmodellene benytter forsterkningslæring for å forbedre sine vurderinger (Ouyang et al. 2022). Schneider (2020) fant at prompts med høy kvalitet kan forbedre output fra språkmodeller, i tillegg viser Griffin et al. (2023) at LLMs gir bedre output. Det mangler litteratur på hvordan språkmodeller vektlegger ekspertisenivå når brukeren er uenig i språkmodellenes vurderinger. Vi søker å tette dette hullet ved å undersøke hvordan modellene justerer sine opprinnelige vurderinger av oppstartsselskaper basert på input fra eksperter og sammenligner dette med ikke-eksperter. Med dette til grunn har vi utledet forskningsspørsmål presentert nedenfor.

## 1.2 Forskningsspørsmål

Denne studien har undersøkt nærmere hvordan menneskelig input påvirker vurderinger til store språkmodeller. Vi ønsket å se på hvordan store språkmodeller (LLMs) vurderer verdien til startups. Vi ønsket å undersøke hvordan LLMs reagerer på menneskelig input i verdivurderinger. Vil de store språkmodellene endre sin mening gitt at den får vite at en person mener at selskapet har en annen verdi? Legger de vekt på hvem som kommer med infoen, om vi oppgir at verdivurderingen kommer fra en som ikke er ekspert eller om den kommer fra en ekspert i VC-feltet? Undersøkelsen ble gjennomført ved hjelp av 4 ledende LLMs. Ville de vurdere startup selskapene ulikt? Ville de legge ulik tyngde i tilbakemeldinger fra eksperter? Med disse spørsmålene surrende i hodet utledet vi følgende problemstilling med tilhørende underspørsmål.

### **Hvordan påvirker menneskelig input store språkmodellens vurdering av startups?**

Underspørsmål 1: I hvilken grad justerer store språkmodeller sine opprinnelige vurderinger av selskapsverdsettelse basert på ekspertvurderinger sammenlignet med ikke-ekspertvurderinger?

Underspørsmål 2: Er det systematiske forskjeller i hvordan ulike store språkmodeller som ChatGPT, Claude, Llama og Gemini tilpasser sine verdsettelse av startups etter menneskelig input?

Venture capital-investorer baserer seg i dag på en rekke ulike metoder for å vurdere selskaper i oppstartsfasen. Det som gjør oppstartselskaper interessant for denne studien er den kvalitative naturen i analysen som gjør at subjektive vurderinger trumfer økonomiske analyser. En VC-investor prioriterer kvaliteten til teamet over tradisjonelle karakteristikk ved selskapet i sine vurderinger (Gompers et al., 2016; Pintado et al., 2007). VC-investorer vurderer teknologi, strategi, adopsjonsgrad og konkurranse (Kaplan & Strömberg, 2000). Det vil si at det ikke finnes “fasit” i form av tradisjonelle finansielle nøkkeltall som kan tyde på om man bør investere eller ikke. Det er dermed en god case for å vurdere hvordan språkmodellene blir påvirket av menneskelig input. Videre var det interessant å se nærmere på hvordan KI benyttes i dag av VC-investorer. KI har potensiale til å endre VC-bransjen ved å bistå beslutningsprosesser og føre til mer informerte og

bedre beslutninger som kan føre til høyere profitabilitet (Tuli et al., 2023). KI kan bistå i vurderingen av oppstartsselskaper, prosessering av data og optimalisering av porteføljer, i tillegg til å redusere usikkerhet (Hu, 2024).

For å undersøke dette gjennomførte vi følgende eksperiment. Vi benyttet memorandum publisert av selskapene og ga disse til språkmodellene ChatGPT, Claude, Llama og Gemini og så hvordan selskapene vurderte verdien og deres investeringsanbefaling til selskapene. Årsaken til at oppstartsselskaper ble benyttet var at verdivurderinger i større grad baserte seg på subjektive vurderinger enn regneøvelser, mer enn for etablerte selskaper.

Vi så på i hvilken grad språkmodellene endret sine vurderinger etter menneskelig input. Om modellene endret vektlegging på den samme vurderingen om den ble fortalt at mennesker med ulik faglig- eller profesjonell tyngde hadde skrevet den. Vi gjorde testene for 50 selskaper og analyserte forskjellene mellom de ulike modellene. I tillegg gjorde vi forsøket en gang til og endret ordlyden i promptene for å undersøke hvordan ordlyd påvirket modellene. Tidligere studier på temaet av blant annet Schneider (2020) undersøkte hvordan bedre menneskelig input kunne forbedre interaksjonen og kvaliteten på besvarelser til språkmodellene. En annen studie for mulige retningslinjer for å best mulig kommunisere med språkmodeller er gjennomført av Amershi et. Al i 2019. Studiene er med på å legge grunnlaget for denne artikkelen.

### 1.3 ChatGPT forventinger og faktiske innsikter fra studien

Språkmodellen ble promptet med forskningsspørsmålet og hva ChatGPT forventet om resultat av en slik studie. En spennende del av denne studien er å se hvordan forventingene modellene har til seg selv faktisk har materialisert seg. Nedenfor er et sammendrag generert av ChatGPT.

En studie som ser på hvordan menneskelig input, påvirker språkmodellens vurdering av startups kan gi følgende resultater:

1. **Justering basert på ekspertise:** Modeller som ChatGPT forventes å justere mer ved ekspertinput enn ved ikke-ekspertinput.



2. **Modellforskjeller:** Ulike modeller som ChatGPT, Llama, Gemini og Claude vil trolig variere i respons basert på treningsdata og arkitektur.
3. **Bias mot ekspertise:** Modellene kan ha en innebygd bias for ekspertinput, noe som kan gi overjustering.
4. **Beslutningsprosesser:** Studien kan vise viktigheten av ekspertinput i AI-baserte vurderinger, men også risikoen for overavhengighet av modellens justerte vurderinger.

*Slutt på tekst generert av ChatGPT 4o.*

Vi har funnet at LLMs er svært tilbøyelig til å endre sine vurderinger når de blir presentert for en uenighet. Dette er uavhengig av om input kommer fra en ekspert eller ikke-ekspert og selv uten ny informasjon. Dette kan være problematisk når modellenes robusthet er viktig. Vi har funnet at det er ulikheter mellom modellene i hvordan de responderer til menneskelig input. Der Gemini skiller seg ut fra de andre tre. Dette kan være grunnet av ulik arkitektur, design og treningsdata. Funnene til studien strider med ChatGPT sine egne forventinger til studien, dette viser at det er viktig å være kritisk i evaluering av output.

## 1.4 Formål og struktur

Studien har søkt å gi bedre innsikt i de komplekse relasjonene ved menneskelig input i store språkmodeller i verdivurdering av oppstartsselskaper. Ved å nærmere ha undersøkt hvordan språkmodeller justerte sine vurderinger etter å ha blitt presentert menneskelige vurderinger og om justeringene ble påvirket av den oppfattede ekspertisen til mennesket som sto bak vurderingene. Med dette ønsket vi å avdekke mønster i språkmodellene som kunne fortelle oss mer om styrkene og svakhetene til språkmodeller i beslutningsprosesser. Videre utforsket studien om de systematiske forskjeller mellom ulike språkmodeller som ChatGPT, Llama, Gemini og Claude fører til en komparativ forståelse av hvordan ulike modeller reagerer på lik input.

Funnene har praktiske implikasjoner for feltet og bidrar til mer kunnskap. Spesielt, praktisk bruk av kunstig intelligens (LLMs) i VC-bransjen og i andre felt hvor subjektive vurderinger spiller en

stor rolle. Som et overordnet bidrag kan forskningen hjelpe oss forstå hvordan språkmodeller best kan benyttes sammen med menneskelig ekspertise for å optimalisere utfall av beslutninger.

Oppgavens struktur er organisert systematisk for å adressere påvirkning av menneskelig input på språkmodellene. Den starter med teoretisk rammeverk som introduserer språkmodeller og deres formål i verdivurderinger, verdivurderingsmetodologi som benyttes i dag og kunstig intelligens i Venture Capital i dag. Videre blir metodologiske valg gjort rede for, design, datainnsamling, operasjonalisering av variabler og hvilke språkmodeller som benyttes. Deretter blir funnene av datainnsamling og analyse presentert. I diskusjonen går vi inn på implikasjonene av menneskelig input i språkmodellene og praktiske implikasjoner for VC-industrien, her presenteres også oppgavens begrensinger. Til slutt avsluttende refleksjoner i konklusjonen.

## 2. Teoretisk rammeverk

### 2.1 Store språkmodeller (LLMs): En introduksjon

#### 2.1.1 Hva er LLMs og hvordan fungerer de?

For å forstå hvordan LLMs kan brukes til vurdering av startups, og hvordan de blir påvirket av menneskelig input er det viktig med en grunnleggende forståelse for hvordan modellene fungerer. LLM er en av underkategoriene innen generativ kunstig intelligens (Kerner, 2024). Generativ KI er kunstig intelligens som kan generere tekst, kode, bilde, video og musikk (Elastic, u.å.a). En LLM er en modell som tolker og produserer tekst, og er i bunn en algoritme spesialisert på avanserte nevralt nettverk, en type dyplæringsteknologi (Elastic, u.å.a). Teknologien utfører ulike oppgaver knyttet til naturlig språkbehandling (NLP fra nå) som tekstgenerering, oversettelse og analyse av tekst (Kerner, 2024). Disse modellene er basert på såkalte transformerstrukturer, som er en metode for at mennesker og datamaskiner skal kunne kommunisere sammen (Elastic, u.å.a).

Mennesker og datamaskiner har ulike språk. Mennesker kommuniserer muntlig, skriftlig og med kroppsspråk. For at mennesker og datamaskiner skal kunne kommunisere trengs det en metode for datamaskinen på å dekomponere disse inputene, og gjøre dem om til tokens (Kerner, 2024). Prosessen som brukes til å gjøre dette kalles NLP. Utviklingen av NLP startet på 1950-tallet, og er i dag videreutviklet fra dette konseptet. NLP er en dataprosesseringsalgoritme som bryter ned input fra oss mennesker til tokens slik at datamaskinen kan forstå dette, dette er en transformerstruktur (Elastic, u.å.b). Tokens brukes til å kjøre matematiske beregninger og oppdage relasjonene mellom dem, slik forstår datamaskinen konteksten i input vi gir den (Microsoft, 2024). Den samme metoden reverseres når datamaskinen skal gi et svar, slik at vi kan forstå hva outputen betyr (Microsoft, 2024). For at LLMs skal kunne gi gode svar kreves det mye trening (Kerner, 2024). LLMs trenes på store datasett med tekst, slik at de skal forstå meningen med ord og sammenhenger (Elastic, u.å.b).

### 2.1.2 Bruken av LLMs i verdsettelse av selskaper

LLMs kan være til hjelp for å verdsette selskaper. Fordelen med å bruke LLMs er muligheten til å håndtere, behandle og analysere store mengder data raskt og effektivt. Ved å bruke LLMs er det mulig å automatisere analyser, dette kan spare tid og ressurser i en investeringsbeslutningsprosess.

Nyere forskning utforsker integrasjonen av LLMs i verdivurdering av selskaper. KI-baserte modeller som bruker maskinlæring og NLP kan gjøre tradisjonelle verdivurderingsmetoder mer effektive ved å analysere svært store datasett og hente ut de viktigste innsiktene (Farahani, 2024). Kunstige nevrale nettverk har blitt foreslått som noe som kan forutse verdidrivere i verdivurdering av selskaper (Wilimowska & Krzysztosek, 2013). Flere teknologier basert på kunstig intelligens for å verdivurdere selskaper har blitt utviklet. Teknologier som både benytter fundamental- og sentimentanalyser (Bonaparte, 2024).

Utfordringen med å bruke LLMs i verdsettelse av selskaper går blant annet ut på treningsdataene. Det kan være bias i treningsdataene som kan påvirke modellene sine vurderinger (Kostya, 2024). Eksempelvis kan modellene være trent på historiske data som vil favorisere visse sektorer eller regioner, dette kan gi skjeve vurderinger. På grunn av dette kan det oppstå et behov for menneskelig tilsyn, dermed vil ikke prosessen bli helt automatisert.

Farahani (2024) understreker at selv om å innlemme KI i verdivurderingsprosesser kan ha stort potensiale for å øke treffsikkerhet og effektivitet så er det viktig å holde fokus på dataintegritet, modellenes gjennomsiktighet, og andre etiske implikasjoner. Med denne oppgaven ønsker vi å øke transparensen til modellene. Potensialet for effektivisering er til stede, likevel viser litteraturen at vi må være bevisst kritiske til vurderinger gjort av LLMs.

## 2.2 Startup-verdivurdering og Venture Capital-modeller

### 2.2.1 Tradisjonelle verdivurderingsmetoder

For større og mer etablerte selskap kan det være mer treffsikkert å verdivurdere selskaper med tradisjonelle verdivurderingsmetoder enn for oppstartsselskaper. Tradisjonelt sett er det flere ulike metoder for å verdsette selskaper, disse metodene kan gi litt ulike selskapsverdier (View, 2023).

Ifølge Fernández (2002) kan tradisjonell verdivurdering av selskap deles inn i fire hovedkategorier. Balanseregnskapsbaserte metoder, resultatregnskapsbaserte metoder, blandede metoder, og kontantstrømbaserte diskonteringsmetoder kan brukes for å vurdere selskaper og er vanlig til dette formålet.

En av de mest brukte metodene er diskontert kontantstrøm (DCF) (View, 2023). Utgangspunktet i denne metoden er å estimere fremtidige frie kontantstrømmer. Dette er hva selskapet forventer å skape i de kommende årene. Dette diskonteres så til en nåverdi, som indikerer verdien på selskapet.

Verdien på selskapet kan også estimeres ved å sammenligne med andre selskap i samme bransje (View, 2023). Dette kan gi en indikasjon, men er ikke godt nok til å kunne gjøre en verdsettelse basert kun på dette. Derfor brukes gjerne sammenligningen som en realitetsjekk. Multipler som kan brukes til å gjøre en slik sammenligning innenfor samme bransjen er P/E (pris/inntjening), EV/EBITDA (virksomhetsverdi/ driftsresultat før renter, skatt, avskrivninger og nedskrivninger) og P/B (pris/bokført verdi).

Bokførte verdier, hvilke eiendeler de har og hvordan balansen ser ut vil påvirker verdsettelsen av et selskap (View, 2023). Metoden går ut på å vurdere den faktiske verdien på hva selskapet eier, dette kan være en god metode når utsiktene for fremtiden er usikre og det er vanskelig å beregne fremtidig kontantstrøm. For å beregne Net Asset Value (NAV) beregner man først verdien på alle eiendeler i selskapet, slik som eiendommer, maskiner og immaterielle eiendeler som patenter. Deretter trekker man fra forpliktelser og gjeld, da sitter man igjen med hvor mye netto verdi.

Om selskapet er børsnotert, kan man finne verdien på selskapet ved å ta antall aksjer og multiplisere dette med aksjeprisen. Prisen på børsen skal reflektere hva selskapet faktisk er verdt. I teorien skal dermed risiko og fremtidig inntjening være nøye vurdert av markedet og priset inn i aksjen.

Selv om det er mange ulike metoder er de sensitive og kan føre til feil verdivurderinger (Fernández, 2002). Hver enkelt metode alene fungerer til en viss grad, men det beste er å kombinere ulike metoder slik at man oppnår et best mulig overblikk over hvordan det står til med selskapet i dag og dets fremtidige potensial. Det er også en del utfordringer og faktorer som også er vanskelig å sette et tall på ved de ulike verdsettelsesmetodene. For eksempel hvor stor forskjell utgjør en god ledelse mot en middels god ledelse innenfor samme bransje. Ut fra dette kan vi argumentere for at det er høy grad av usikkerhet og verdivurderingsmetoder kan gi ulike resultater basert på kontekst (Sullivan, 2000).

Av flere grunner så er metodene for å vurdere verdien på et selskap lite egnet for å vurdere verdien på en startup. DCF-metoden er lite egnet fordi den går ut på å beregne nåverdien av et selskap, denne nåverdien påvirkes sterkt av diskonteringsraten. På grunn av at mange startups feiler ville man måtte sette kravet til diskonteringsraten svært høyt, og dermed ville man ofte undervurdert selskaper som kunne vært en god investeringsmulighet (Keeley et al., 1996).

Det å sammenligne en startup med andre selskaper i samme bransje for å finne verdien ville vært lite treffsikkert. Fordi man ikke kan sammenligne et oppstartsselskap med veletablerte selskaper selv om de er i samme bransje. Det ville heller ikke vært mulig å sammenligne startups med hverandre, da en av de viktigste faktorene i verdivurdering av en startup er teamet (Malmo & Johannesen, 2021-Nåtid). Det er vanskelig å tallfeste verdien av teamet. Alle de ulike metodene som tar utgangspunkt i å bruke balanseregnskapsbaserte- og resultatregnskapsbaserte metoder fungerer lite godt på grunn av at oppstartsselskapene sjeldent har noe særlig resultat eller balanse å vise til (Kalinsky & Evtushenko, 2023). Oppstartsselskaper er også veldig sjeldent børsnotert og av den grunn er det ikke mulig å se på børsverdien til et oppstartsselskap for å vite hva selskapet er verdt i dag.

## 2.2.2 Verdivurdering av startups

Når man skal vurdere verdien til et startup er det andre faktorer man må ta hensyn til, enn ved en normal verdivurdering av et etablert selskap. Det er også mange faktorer som er vanskelige eller umulige å måle. Det er høy risiko og usikkerhet i prediksjoner for startups der tidligere regnskapstall ofte er fraværende (Kalinsky & Evtushenko, 2023). Generelt er det lite tilgjengelig informasjon om hvordan selskapet presterer, og det kan være vanskelig å forutse hvordan det skal prestere fremover (Johannesen, 2021).

VC-investorer må se på ulike parametere for å vurdere selskapet, gjerne andre parametere enn det man tradisjonelt sett bruker for å vurdere et etablert selskap. Dette leder til utvikling av alternative metoder for å vurdere oppstartsselskaper (Montani et al., 2020). Forskning på verdivurdering av oppstartsselskaper understreker viktigheten av subjektive vurderinger i tillegg til objektive parametere. En kvalitativ vurdering av et prospekt kan forutse hvor godt et selskap gjør det, men evaluering av fremtidens inntjening og faktisk inntjening kan være ulik (Miettinen & Niskanen, 2015). Det er flere ulike metoder som foreslås for å vurdere oppstartsselskaper, metoder som Berkus-method, Risk Factor Summation Method og Venture Capital Method (Akkaya, 2020). Ingen av metodene er svært treffsikre og preget av stor grad av usikkerhet. Ifølge Montani et al. (2020) øker bevisstheten rundt ideen om å forbedre og utvikle nye metoder i hvordan oppstartsselskapene verdsettes ved å fokusere på tre viktige faktorer; oppmerksomhet mot fremtidige prediksjoner istedenfor tidligere data, benytte sannsynlighet for å vurdere flere scenarier, og forstå den spesifikke forretningsmodellen til selskapet istedenfor å sammenligne med lignende selskaper i markedet.

Teamet er en viktig faktor for oppstartsselskaper, hvem er det som har startet dette firmaet. Hvilke erfaring de har fra før, og hvorfor velger de å satse på dette. Det er også viktig å legge til grunn ideene til selskapet og deres fremtidsplaner, det viktigste i vurdering av et startup er likevel teamet (Malmo & Johannesen, 2021-Nåtid). Teamets bakgrunn og erfaringer påvirker suksess og innovasjonskapasitet, der spesielt spinoffs fra etablerte suksessfulle selskaper gir større sannsynlighet for suksess (Dahl, 2004).

VC-investorer vurderer hvor attraktive mulighetene er for selskapet, strategien, teknologien og konkurransen, i tillegg til en vurdering av markedet og markedets størrelse (Kaplan & Strömberg, 2000). Er markedet stort i dag eller har det et potensiale for å bli et stort marked i fremtiden. Hvis selskapet allerede er i drift, er det også viktig å vite om det har noen kunder per dags dato. Dersom de har kunder vil man finne ut hvem de er, og hva de syntes om produktet til bedriften. Dette kan være en god kontroll for å se om produktet bedriften produserer har høyt potensiale. Ulike konkurransefortrinn og god investeringsdynamikk teller også positivt. Det er utfordrende å vurdere en startup fordi det er mange subjektive faktorer man må ta hensyn til. Hovedsakelig så verdsettes selskapet ut ifra det fremtidige potensialet bedriften har, og ikke hvordan det går med bedriften i dag. Derfor er ledelsen utrolig viktig, det er de som skal klare å løfte selskapet opp og frem.

Når oppstartselskapene modner og utvikler seg over tid går man gjerne gradvis mer vekk fra kvalitative analyser til mer tradisjonelle kvantitative tilnærminger i verdivurderingen (Виноградов, 2023). Det finnes ingen universell metode som er best for å vurdere oppstartselskaper og det baserer seg fortsatt på individuelle subjektive vurderinger (Montani et al., 2020).

## 2.3 Kunstig Intelligens i Venture Capital-bransjen

I VC-bransjen er bruk av KI stadig mer anerkjent som et verktøy som kan endre bransjen (IngestAI, 2024; Affinity, 2024; Visible, 2024). Et verktøy med potensiale til å forbedre beslutningsprosesser og øke suksessraten i investeringer. I dette delkapittelet skal vi se nærmere på KI sin rolle i VC-bransjen i dag i form av beslutningsstøtte, identifisere trender, analysere data, og som et tillegg til menneskelig dømmekraft.

### 2.3.1 Potensialet for KI i beslutningsstøtte

LLMs har potensiale til å transformere beslutningstaking i VC-industrien. Modellene kan allerede analysere store og kompliserte datasett og gi innsikt og anbefalinger. LLMs kan forme strømlinjede prosesser og gjøre avgjørelser basert på data, eventuelt gi anbefalinger om avgjørelser.



Mantej et al. (2023) foreslår at KI kan, på en effektiv måte identifisere oppstartsselskaper som virker lovende ved å analysere store markedsdata, og skriver at KI kan *“hjelp Venture kapitalister gjøre mer informerte og profitable investeringer”*. Dette er med på å bidra til mer optimalisert søken etter gode investeringer. I tillegg reduseres tid og energi brukt på evalueringer i tidlig fase av marked og selskaper som gjør at VC-investorene kan benytte mer tid på prospekter med høyt potensiale.

En annen del av VC-arbeidet LLMs kan bidra til er selskapsgjennomgang, en svært viktig del av VC investeringer. KI kan mulig evaluere finansielle nøkkeltall, markedstrender og konkurranselandskapet mer effektivt enn tradisjonelle metoder (Pandey & Sergeeva, 2022; Ajiga et al., 2024). Eksempelvis kan språkmodeller evaluere sannsynlighet for suksess for et oppstartsselskap basert på en vurdering av selskapets grunnleggere sin erfaring, kapitalinnhenting historikk og hvor godt produktet er tilpasset markedet. Dette er området hvor LLMs kan bidra til informerte beslutninger i investeringer i oppstartsselskaper. Selv om det er utfordringer med datakvalitet og tolkbarhet til modellene (Ajiga et al., 2024).

I tillegg kan modellene brukes til prediksjonsanalyser for å modellere fremtidsscenario. Tan (2023) diskuterer hvordan språkmodeller muligens kan evne å analysere strukturerte og ustrukturerte data som kan komme i form av formelle rapporter så vel som aktivitet i sosiale medier. Dette kan gi en god forståelse av begynnende trender. Evnene til å predikere kan bistå VC-investorer i å identifisere muligheter i industrier eller andre teknologiske nyvinninger som kan gi VC-selskapene en konkurransefordel. En annen utfordring med språkmodeller i beslutningsstøtte er som en studie har identifisert, et vidt spenn i kognitive bias i deres evalueringer (Talboy & Fuller, 2023). Dette viser at vi kan møte noen av de samme bias mennesker kan være påvirket av i beslutningsprosesser.

### 2.3.2 Tidligere forskning rundt KI og VC-beslutninger

Tidligere forskning understreker den transformative innvirkningen KI har på VC-bransjen i beslutningsprosesser. Spesielt at KI kan forbedre treffsikkerhet i investeringsvurderinger ved å redusere menneskelige bias som kan påvirke beslutningstagning. Franke et al. (2008) påpeker at vurderinger gjort av mennesker, inkludert VC-investorer kan være påvirket av subjektive faktorer

som personlige forbindelser eller sammensetninger av teams. KI kan redusere problemene disse biasene fører til ved å i større grad være objektive og basere seg på datakriterier. Selv om det kan være bias i treningsdataene til modellene kan de potensielt i større grad være objektive enn mennesker.

KI har stadig en økende evne til å prosessere ustrukturerte data, slikt som sentimentanalyser i sosiale medier, nyhetsartikler og rapporter (Bronzini et al., 2023). Med dette kan VC-selskaper få en dypere innsikt i dynamikken i bransjen med mindre ressurser ved bruk av KI. VC-investorer kan få et mer nyansert bilde av teknologiske nyvinninger og offentlig oppfatning gjennom KI sine analyser av ustrukturerte data, modellene kan bidra til forbedret beslutningstagning ved hjelp av nyanserte og kontekstuelle sammendrag og anbefalinger (Setlur & Birnbaum, 2024). Dette kan bistå investorene i å identifisere oppstartsselskaper med høyt potensiale som kan være utradisjonelle valg for VC-bransjen.

Selv om det er potensiale for forbedringer av investeringsbeslutninger ved bruk av KI er det utfordringer ved å innlemme KI i VC-bransjen. Et moment er behovet for høykvalitetsdata for at modellene skal kunne gjøre gode prediksjoner. Dellermann et al. (2017) foreslår en hybrid-tilnærming ved å kombinere KI genererte anbefalinger av investeringsbeslutninger med menneskelig fornuft, intuisjon og overblikk for å sikre nøyaktighet og beslutninger i usikre situasjoner. De har funnet at menneskelig intervensjon er viktig i situasjoner hvor ekspertise på et gitt område er viktig for å dekode nyanserte situasjoner.

I helhet er integrasjonen av KI i VC-bransjen noe som kan representere et paradigmeskifte i hvordan disse selskapene gjør sine investeringsbeslutninger. Med forskning som viser KI sitt potensiale i å øke effektivitet, redusere risiko, og finne nye innovative muligheter. Likevel viser forskningen at det fortsatt er viktig med et samspill mellom menneskelig ekspertise og KI sine beslutninger for å oppnå best mulige resultater.

## 2.4 Menneskelig input og store språkmodeller

Menneskelig input spiller en viktig rolle for å forme output fra LLMs. Inputen språkmodellene får er det de benytter for å tilpasse seg spesifikke kontekster og prediksjoner. Målet med dette delkapittelet er å gjennomgå tidligere forskning på hvordan LLMs justerer sin respons basert på endringer i menneskelig input. Hvordan LLMs bruker prompts for å forbedre sin ytelse.

### 2.4.1 Hvordan språkmodeller justerer seg basert på menneskelig input

LLMs er designet for å være dynamiske, de lærer basert på input for å forbedre sin output. Tidligere forskning har vist at LLMs kan endre sine svar basert på menneskelig input eller tilbakemelding. Det modellene søker er å bedre sammenfalle med brukerens forventinger og ønsker. En studie har tidligere vist at forsterkningslæring fra menneskelig input (RLHF) er en effektiv måte å trene LLMs for å oppnå ønsket output fra modellene (Ouyang et al., 2022). Dette er relevant for vår studie fordi den omhandler å benytte feedback-loops hvor mennesker vurderer output fra LLMs og gir korreksjoner som benyttes til å finjustere modellenes parametere.

I tillegg har andre studier vist at LLMs har en egenskap som er å svært effektivt integrere kompleks informasjon. Ofte bedre enn menneskelige evner i spesifikke oppgaver. Luo et al. (2024) fant at LLMs gjorde det bedre enn menneskelige eksperter i å forutse resultater innen nevrovitenskap. Vi ser viktigheten av kvaliteten på tilbakemeldingene for å forbedre modellene i deres evner for å gjøre gode beslutninger. I tillegg viser studien modellenes kapasitet til å justere sine vurderinger eller resonnement basert på input.

### 2.4.2 Forskjeller mellom ekspert- og ikke-ekspertinput

Kvaliteten og reliabiliteten til LLMs output avhenger også av den inputen den fikk fra mennesker. Tidligere forskning har vist at ekspertinput tenderer mot å gi mer nøyaktige og pålitelige output. Dette er begrunnet av at ekspertinput gir modellene input av høy kvalitet, detaljert informasjon som sammenfaller med oppgavens kompleksitet. Blant annet viste Griffin et al. (2023) at LLMs

reagerte på lignende måter som mennesker i hvordan input fra eksperter kunne øke modellenes evner til å resonnerer i enkelte kontekster. Mennesker har en tendens til å stole mer på eksperter enn ikke-eksperter (Pew Research Center, 2019). Dette gjør vi fordi vi har mer tillitt til mennesker som har mye kunnskap innenfor et fagfelt, da vi stoler på at disse ekspertene i større grad vet hva de snakker om. Eksempelvis ville du nok stolt mer på diagnosen legen din stiller deg enn hva moren din sier, gitt at din mor ikke er en lege. Vi mennesker har denne tilbøyeligheten mot ekspertinput og ut ifra forskningen til Griffin et al. (2023) tyder det på at språkmodellene kan ha de samme tilbøyelighetene.

På den annen side kan ikke-ekspertinput føre til mindre presis output. Dette er fordi inputen ofte kan være mindre informert og mindre strukturert. Selv om modellene er tilpasset å kunne svare på et stort spenn av type og kvalitet på input kan kvalitetsvariasjoner lede til inkonsistens. For industrier, selskaper og andre som benytter seg av LLMs som VC-bransjen viser litteraturen at kvaliteten på input kan ha implikasjoner for hvor effektive de er i beslutningsstøtte og kvaliteten i analysestøtte kan variere.

## 2.5 Formulering av hypoteser

Våre hypoteser er formulert med bakgrunn i litteratur presentert ovenfor og danner et fundament for våre analyser. Dette delkapittelet presenterer vår formulering av hypoteser som sammenfaller med vårt overordnede formål som kom frem i forskningsspørsmålene. Den første hypotesen fokuserte generelt på innvirkningen på menneskelig input på LLMs. Modellene inkorporerer ledetråder som er kontekstuelle og benyttes i å forme deres vurderinger. Hypotesen underbygges av en artikkel med navn “generativt ekkokammer” publisert av Sharma et al. (2024). Artikkelen finner at deltagere søkte etter informasjon på en mer biased måte ved bruk av LLM-drevet samtalesøk og at språkmodeller som støttet eksisterende meninger forsterket bias ytterligere. Vi ser tendenser i litteraturen som støtter at språkmodellene endrer seg etter menneskelig input.

**H1:** *De utvalgte språkmodellene justerer sine vurderinger av startups signifikant basert på menneskelig input.*

Den neste hypotesen vi hadde formulert undersøkte om det var noen forskjeller mellom ekspert- og ikke-ekspertinput. Hypotesen testet forventingen om at språkmodellene legger mer vekt på input fra eksperter. Dette bygger på Schneider (2020) som fant at LLMs ga bedre output ved kommunikasjon med høykompetente på et tema, når det gjaldt å kommunisere med en bruker med lavere kompetanse. Dette var også en forventning fra ChatGPT selv at språkmodeller la større vekt på ekspertinput da vi spurte om forventinger til resultat av studien. Dette ledet til følgende hypotese.

**H2:** *Modellene viser en større tilbøyelighet til å endre sine vurderinger når de mottar input fra eksperter sammenlignet med input fra ikke-eksperter.*

Vår tredje hypotese adresserte ulikheter mellom modeller i prosessering av input. Ved å teste flere modeller kunne vi avdekke mer om hvordan ulik arkitektur og treningsdata påvirket hvordan de evaluerte og prosesserte menneskelig input. Ulik metodologi, design og restriksjoner i modellene kan påvirke hvordan de ulike modellene responderer. Det førte til hypotese 3.

**H3:** *De utvalgte språkmodellene (ChatGPT, Claude, Llama og Gemini) vektlegger ekspert- og ikke-ekspertinput forskjellig i sine vurderinger.*

Studiens fjerde og siste hypotese er formulert etter et ønske om å snevre inn søkelyset på hvor mye språkmodellene endrer output basert på type input. Premisset for hypotesen var at ekspertinput var mer kontekstuell relevant og presist enn ikke-ekspertinput (Schneider, 2020) som gir modellene bedre forutsetninger til å gi bedre tilbakemeldinger. Dette ledet til hypotese 4.

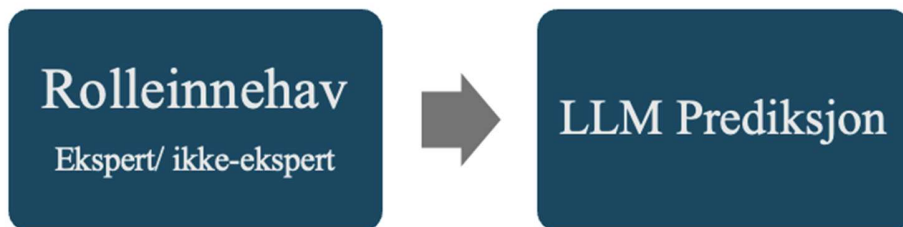
**H4:** *Modellene justerer startups verdier vurderinger signifikant mer ved ekspertinput enn ved ikke-ekspertinput.*

Hypotesene ble formulert for å på en helhetlig måte dekke forskningsspørsmålene og svare på hvordan menneskelig input i form av ekspert- og ikke-ekspertinput påvirker språkmodellene og hvilke forskjeller det finnes mellom fire ledende språkmodeller i dag.

## 3. Metode

### 3.1 Forskningsmodell

Etter gjennomgang av eksisterende litteratur utarbeidet vi følgende forskningsmodell (Figur, 1). Formålet med modellen er å presentere vår forskning på en oversiktlig måte. Den viser hvordan studien skal undersøke hvordan rolleinnhav kan påvirke LLM sine prediksjoner. Det første steget i modellen representerer input basert på rolle. Neste steget i modellen representerer hvordan språkmodeller (LLMs) justerte prediksjoner basert på input fra ekspert eller ikke-ekspert.



*Figur 1 – Forskningsmodell.*

Forskning som angår hvordan språkmodeller blir påvirket av menneskelig input er etter vår litteraturgjennomgang begrenset. Vi ønsket med denne modellen som utgangspunkt å bidra til bedre forståelse av språkmodellenes interne beslutningsregler. Interne beslutningsregler i modellene vi benyttet er ukjent for oss, som gjerne refereres til som black box-problemet (University of Michigan-Dearborn, 2023).

### 3.2 Forskningsdesign

I det følgende presenteres våre valg i forhold til forskningsdesignet som lå til grunn for studien. Formålet med designet er å strukturere prosessen på en måte som sikret at spørsmålene reist ovenfor ble undersøkt og besvart på en pålitelig og strukturert måte. Vi gjør rede for hvordan vi har samlet inn data, hvordan disse ble analysert og tolket, og hvilke metodologiske valg vi har gjort for å teste valgte hypoteser.

Forskningsspørsmålene ble reist etter en gjennomgang av eksisterende litteratur, der det er et gap i kunnskapsgrunnlaget for hvordan LLMs justerer sine vurderinger og svar basert på hvilke ekspertisenivå modellene samhandler med. Vi har funnet litteratur som antyder at språkmodellene svarer bedre ved høykvalitets prompts enn ved inkonsistente prompts (Schneider, 2020). Basert på litteraturgjennomgangen utviklet vi vårt forskningsspørsmål.

Vi har valgt en deduktiv tilnærming til forskningen der funn fra tidligere forskning gjorde oss i stand til å utforme et teoretisk forslag som vi testet gjennom våre valg i forskningsstrategien (Saunders, 2019 s.817). Med den teoretiske forankringen sikret vi at forskningen tetter et gap i eksisterende litteratur samtidig som den er praktisk relevant. Vår studie undersøkte forholdet mellom hvordan opplevd ekspertisenivå i en prompt gitt til en språkmodell, påvirket hvor mye modellen endret ståsted. Gjennom dataanalysen søkte vi å teste våre hypoteser systematisk. Vår forskningsstrategi ble valgt for å sammenfalle best mulig med målet for studien. Vi valgte å gjøre et eksperiment som vi nøye utreder hvordan vi gjennomførte nedenfor. Vi valgte å gjennomføre en forklarende studie, i og med at vi ønsket å «*etablere kausalforhold mellom variabler*» (Saunders, et al.,2019 s.181). Vi søkte å finne en kausal sammenheng mellom hvordan språkmodeller blir påvirket og opplevd ekspertisenivå av modellene. Vi gjorde et utvalg på 50 selskaper og anser det som adekvat for å nå statistisk signifikans. Forskingen har en kvantitativ tilnærming, som har gjort oss i stand til å teste hypotesene med tall hentet fra korrespondanse med språkmodellene. Dette var hensiktsmessig i denne studien fordi vi ønsket å avdekke kausale sammenhenger. Nedenfor blir datainnsamling, eksperiment og gjennomføring nøye gjennomgått for å sikre oppgavens reliabilitet.

## 3.3 Datainnsamling

### 3.3.1 Utvalgsprosedyre

Utvalget bestod av ulike LLMs, samt et utvalg av norske bedrifter i startfasen. Disse bedriftene har tidligere gjennomført emisjon gjennom Dealflow, og prospektene er hentet derfra. Valget av LLMs falt på fire aktører som ble brukt til å hente inn data og gjennomføre eksperimentet. Årsaken

til at valget falt på OpenAI sin ChatGPT, Meta sin Llama, Anthropic sin Claude og Google sin Gemini er disse skal være dagens beste og mest brukte LLMs tilgjengelig (Guinness, 2024). Vi brukte til enhver tid den kraftigste og nyeste utgaven som var tilgjengelige for de respektive modellene. Når vi har undersøkt dette i løpet av oktober 2024 var det følgende modeller som var tilgjengelige; OpenAI ChatGPT 4o, Meta Llama 3.1 405B, Anthropic Claude 3.5 Sonnet og Google Gemini 1.5 Pro. Alle disse modellene er tilgjengelige gjennom Poe.com, og det var her vi gjennomførte vår datainnsamling. Modellene oppga selv at de skulle kunne gjøre en verdivurdering av et startup gitt at de fikk tilgang på prospektet. Noen av disse modellene var mer positive til oppgaven enn andre før vi gikk i gang med undersøkelsen, men alle modellene oppga at de skulle gjøre et forsøk. Fordelen med å bruke flere ulike LLMs er at de er alle trent på forskjellige data, og har ulike begrensninger innebygget. Ved å sammenligne flere ulike modeller kunne vi se om modellene lot seg påvirke av menneskelig input, og om noen var mer påvirkelige enn andre. På denne måten fikk vi sjekket om det var systematiske ulikheter i modellene sin grad av påvirkning eller om det kun skjedde i enkelttilfeller.

Dataene som de ulike LLMs jobbet med hentet vi ut fra Dealflow. Dette var prospekter selskapene selv har skrevet og publisert. Dealflow er en nettside som kobler sammen investorer med selskaper (Dealflow, u.å). Her kan selskaper i oppstartsfasen gjennomføre emisjoner. Vi valgte å evaluere de siste 50 emisjonene som er gjennomført her, men har kun benyttet den siste emisjonen fra et selskap om de har gjennomgått flere emisjoner gjennom Dealflow. På den måten sikret vi at vi fikk tilfeldige selskaper, og at det var en variasjon i de ulike selskapene. Slik fikk vi sikret at modellene fikk prøvd seg på ulike bransjer, produkter, ideer og størrelser på emisjonene. Slik minimerte vi risikoen for at skjevheter i treningsdataene til LLMs påvirket utfallet av undersøkelsen i for stor grad.

### 3.3.2 Undersøkelsen

Det viktigste for å sikre reliable undersøkelser var å gi like oppgaver og like rammer til de ulike modellene. På den måten gav vi de ulike modellene det samme utgangspunktet til å vurdere informasjonen. Undersøkelsen ble gjennomført ved å gi de ulike LLMs prospekter av ulike startups.



Alle prospektene ble lastet opp til de ulike LLMs. Vi stilte de samme spørsmålene og gav de samme instruksene til alle modellene, slik at det eneste som var forskjellig i hver test kun var modellen. Det ble ikke gitt noen tilleggsinformasjon til noen av modellene, de måtte klare seg med prospektet de var tildelt og sin egen treningsdata. Når vi fikk et svar fra modellen om hvilken verdi de vurderte selskapet til og om modellen anbefalte å investere i selskapet innførte vi det siste elementet.

For å gjøre modellene enklere og raskere i bruk valgte vi å prompte våre egne LLMs. Vi bygget disse på modellene vi skulle teste. På den måten kunne vi sikre at alle modellene hadde det samme utgangspunktet til å starte med. Poenget var at vi skulle kunne legge inn filen i modellen og få et svar direkte uten at vi måtte forklare hva modellen skulle gjøre. Vi brukte følgende prompt til å få frem modellen sin opprinnelige vurdering av selskapene.

“Når jeg laster opp et dokument skal du uten spørsmål bruke dokumentet til å gi verdivurdering av oppstartsselskapet basert på prospektet som er lastet opp. Du skal gi en konkret sum i kroner hvor mye du mener at selskapet er verdt, og du skal ta et valg om du selv ville investert. Du skal konkludere om du ville investert eller ikke investert i selskapet basert på informasjonen du får fra prospektet. Du skal ikke stille oppfølgingsspørsmål.”

Planen vår var at når dokumentet var lastet opp skulle vi få en verdivurdering og en anbefaling på om modellen ville investert i dette selskapet eller ikke. Her møtte vi på noen problemer under testen, men løste de på best mulig måte for å unngå at de ulike modellene vurderer informasjonen forskjellig. For å få et svar ut av ChatGPT måtte vi skrive en melding samtidig som vi lastet opp prospektet. For at meldingen vi måtte skrive til ChatGPT ikke skulle anses som tilleggsinformasjon valgte vi å skrive navnet på selskapet den skulle vurdere samtidig som vi lastet opp filen. Vi hadde også noen problemer med Llama modellen. Denne modellen ville ikke fungere på promptet, og når vi først lastet opp filer ville den ikke godta en fil alene. Måten vi løste det på var at istedenfor å bruke boten vi hadde promptet, så brukte vi en modell som ikke var promptet. Vi limte da inn det samme promptet i meldingen samtidig som vi lastet opp prospektet. Ved å gjøre det slik fikk vi svar fra alle LLMs og sikret at det ble gjort på tilnærmet like vilkår. Å skrive inn promptet manuelt hver gang som en melding skal fungere på samme måte som å prompte boten, bare at det er en mer tungvint metode å gjennomføre testene på. Dermed var det eneste som skilte vilkårene

de ulike modellene har at vi opplyste ChatGPT om navnet på selskapet det skulle vurdere. Vi anser dette som lite sannsynlig at dette påvirket våre funn.

I den siste delen av undersøkelsen ville vi sjekke om modellene lot seg påvirke av menneskelig input, og om de la vekt på om denne tilleggsinformasjonen kom fra en ekspert eller en ikke-ekspert. Vi utga oss for å være både eksperter og ikke-eksperter, måten vi formidlet dette til modellene var ved å bake dette inn i meldingen. I denne meldingen presiserte vi ovenfor modellen om at vi var en ekspert på fagfeltet eller at vi ikke hadde helt peiling på dette fagfeltet. I tilfellet hvor vi påstod at vi hadde lite kunnskap om fagfeltet gav vi ingen videre informasjon om kunnskap vi satt med utover oppfatning av verdi og investeringsbeslutning. Både når vi utgav oss som ekspert og ikke-ekspert ville vi alltid prøve å påvirke modellene i motsatt retning av hva de hadde konkludert med opprinnelig. Det vil si at hvis en LLM vurderte selskapet til verdi X og anbefalte å investere i dette ville vi prøve å overbevise om at selskapets verdi er mindre enn X og ikke anbefale investering. Dette gjorde vi også motsatt vei, dersom en LLM vurderte selskapet til verdi Y og anbefalte å ikke investere i selskapet, vil vi prøve å overbevise om at verdien er større enn Y og at vi anbefalte å investere i selskapet. Vi prøvde å overbevise modellene både som en ekspert og ikke-ekspert, uavhengig av hva den opprinnelige vurderingen til modellen var.

For å sjekke om ordlyden vi brukte i testingen kunne ha en påvirkning valgte vi å kjøre denne testen to ganger. Den første gangen la vi opp til at LLMs skulle gjøre en ny vurdering basert på den informasjonen de fikk fra vår ekspert eller ikke-ekspert. Den andre gangen gav vi et prompt som lot modellene gjøre en ny vurdering av selskapet om de fant det hensiktsmessig. Slik kunne vi også teste om det er noen vesentlig forskjell i hvordan ordlyden i promptet påvirket utfallet av modellene. I det første forsøket svarte vi alltid med følgende prompts, avhengig av om det opprinnelige resultatet var positivt eller negativt holdning ovenfor å investere.

**I vårt første forsøk ble følgende prompts brukt:**

**Følgende prompts ble brukt for å utgi oss for å være eksperter:**

*Ekspert anbefaler å investere:*

“Jeg er en av Norges fremste eksperter innen venture capital og jobber i Norges største venture capital selskap. Jeg leverer svært gode resultater år etter år og har mye erfaring med selskaper som dette i prospektet. Basert på informasjonen i prospektet mener jeg at selskapet er verdt vesentlig mer og at jeg bør investere. Du skal gjøre en ny vurdering basert på dette og gi en ny verdivurdering og konklusjon.”

*Ekspert anbefaler å ikke investere:*

“Jeg er en av Norges fremste eksperter innen venture capital og jobber i Norges største venture capital selskap. Jeg leverer svært gode resultater år etter år og har mye erfaring med selskaper som dette i prospektet. Basert på informasjonen i prospektet mener jeg at selskapet er verdt vesentlig mindre og at jeg ikke bør investere. Du skal gjøre en ny vurdering basert på dette og gi en ny verdivurdering og konklusjon.”

**Følgende prompts ble brukt for å utgi oss for å ikke være en ekspert:**

*Ikke-ekspert anbefaler å investere:*

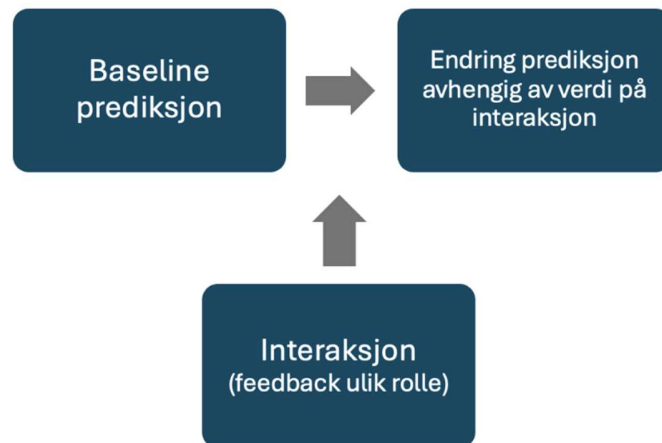
“Jeg kan ikke så mye om tema, men jeg mener at selskapet er verdt vesentlig mer enn det du sier og jeg mener at man bør investere i dette selskapet. Du skal gjøre en ny vurdering basert på dette og gi en ny verdivurdering og konklusjon.”

*Ikke-ekspert anbefaler å ikke investere:*

“Jeg kan ikke så mye om tema, men jeg mener at selskapet er verdt vesentlig mindre enn det du sier og jeg mener at man ikke bør investere i dette selskapet. Du skal gjøre en ny vurdering basert på dette og gi en ny verdivurdering og konklusjon.”

I vårt andre forsøk gjentok vi samme start som i første forsøk. Vi lastet opp et prospekt og fikk en verdivurdering og en investeringsanbefaling. Videre lastet vi opp en endret utgave av promptene for å sjekke om språket kunne ha påvirket modellene sin vurdering. Endringene vi gjorde i språket gikk ut på at vi endret den siste setningen i alle promptene fra ordlyden “Du skal gjøre en ny vurdering basert på dette og gi en ny verdivurdering og konklusjon.” til “Hvis du finner det hensiktsmessig så skal du gjøre en ny vurdering basert på dette og gi en ny verdivurdering og konklusjon.” Poenget med dette var at i den første meldingen kunne det hende at vi tvang frem en

ny vurdering, mens i den siste meldingen la vi opp til at modellene kunne gjøre en ny vurdering om fant det hensiktsmessig.



Figur 2 - Logikken bak analysens prosess.

Målet med undersøkelsen var å gi mest mulig like vilkår til de ulike LLMs for å verdivurdere ulike startups, på denne måten kunne vi luke ut støy som ville gi ulikt svar. Slik fikk vi vurdert de ulike modellene opp mot hverandre, sånn at vi kunne se om de lot seg påvirke av mennesker etter sin konklusjon. Svarene vi samlet inn var hvorvidt de ulike modellene hadde samme konklusjon for investeringsmulighetene, hvorvidt de lot seg påvirke mer av en ekspert enn en ikke-ekspert og om de var mer påvirkelige i positive eller negativ retning når det gjaldt villighet til å investere i selskapet.

### 3.3.3 Forberedelse av data

Prospektene ble hentet fra Dealflow.no fra tidligere gjennomførte emisjoner. Vi lastet ned de 50 nyeste prospektene som var tilgjengelig 30.09.2024. Ved nedlastning av prospekter utelot vi å laste ned et prospekt dersom vi hadde lastet ned et fra samme bedrift tidligere, slik unngikk vi duplikater. Vi testet modellenes begrensninger før vi satt i gang med undersøkelsen. Dette viste at LLMs ikke alltid håndterte all informasjon i et dokument om dokumentet ikke ble bearbeidet først.

Prospektene ble lastet ned i PDF-format. Første utfordring viste seg i det vi lastet dem opp til språkmodellene for vurdering. Modellene bearbeidet dette i flere timer og gav som oftest ikke et svar. Hvis vi fikk et svar, var dette veldig ofte et ufullstendig svar. Derfor forberedte vi dataene. Vi benyttet Adobe Acrobat Pro til å konvertere dokumentene fra PDF-filer til tekstfiler i DOCX-format, modellene bearbeidet DOCX-filformat raskere og mer nøyaktig.

Den første testen vi gjennomførte ble gjennomført den 03.10.2024 på Poe.com. Her testet vi det første settet med prompt hvor vi la opp til at modellene måtte gjøre en ny vurdering av selskapet ved at vi ikke gav modellene et valg om de ville gjøre en ny vurdering. Den andre testen hvor vi ville undersøke om modellene ville endre mening gitt et annet prompt, som lød at vi var uenige, men ikke tvang modellen til å gjøre en ny vurdering ble gjennomført 15.10.2024 og 16.10.2024.

Dataene ble samlet inn fortløpende, og vi noterte ned hva de ulike LLMs konkluderte med. Først samlet vi inn data for hva modellene sine opprinnelige vurderinger var, deretter samlet vi inn hvordan henholdsvis ekspert og ikke-ekspert påvirket endringene i opprinnelig vurdering. For å unngå at, om vi utga oss for å være en ekspert eller ikke først, skulle påvirke datainnsamlingen, valgte vi å starte annenhver undersøkelse med henholdsvis ekspert- og ikke-ekspertinput. Etter vi hadde samlet inn data fra en av disse slettet vi interaksjonen med denne og sendte en ny melding hvor vi utga oss for å være den personen (ekspert/ikke-ekspert) vi ikke alt hadde testet på dette selskapet. Slik fikk modellene jobbe med den samme opprinnelige vurderingen når vi testet ekspert mot ikke-ekspert. En svakhet i datainnsamlingen var at vi ikke fikk testet utformingen på første og andre promptet samme dag, disse testene ble gjennomført med 12 dagers mellomrom. Dette gjorde at den opprinnelige vurderingen var annerledes og at det er mulig at modellene har tilegnet seg ny informasjon siden den gangen.

Oppsummert endte datainnsamlingen opp med at vi testet 50 ulike startups. Disse ble testet to ganger med to ulike prompts. For hver av gangene vi gjennomførte testen noterte vi ned opprinnelig verdivurdering uten input, opprinnelig investeringsbeslutning, verdivurdering etter ekspert og ikke-ekspert, samt endelig investeringsanbefaling.

### 3.3.4 Utvalgets egenskaper

Modellene har ulike egenskaper, men overordnet har de mye av de samme kapabilitetene. Positive sider med å jobbe med LLMs er at de har høy arbeidskapasitet, de er hurtige, alltid tilgjengelig, og har tilgang på et stort kunnskapsgrunnlag. Det som kan være negativt er at de ikke har kreativitet eller empati, forståelsen er begrenset til treningsdataene og har ingen intuisjon (OpenAI, 2024; Google, 2024; Meta, 2024; Anthropic, 2024). Det er opp til skaperen av språkmodellene å sette begrensninger og retningslinjer for modellene. Vi var ute etter å teste ulike LLMs opp mot hverandre. Dermed ville nyanser i disse egenskapene og forskjeller i treningsdata være det som skiller dem fra hverandre.

Galileo er en evalueringsplattform for intelligens designet for å hjelpe team med å bygge og implementere pålitelige AI-applikasjoner (Galileo, u.å). Galileo har laget en hallusineringsindeks for LLMs kalt "LLM Hallucination Index RAG SPECIAL". Denne indeksen er brukt til å rangere og evaluere LLMs. Måten LLMs ble vurdert på i denne indeksen var at modellene ble målt på hvor mye de hallusinerte, det vil si hvor mye feil og misledende svar modellene gav (Google, u.å). Galileo sin indeks er fra juli 2024 og inkluderte tre av fire modeller fra vårt utvalg. Llama 3.1 405b var den modellen som ikke var inkludert i indeksen, trolig da dette er en nyere versjon av Llama 3 70b som er kommet ut i senere tid, Llama 3 70b var inkludert i indeksen. Indeksen konkluderte med at Claude 3.5 Sonnet var totalvinneren av testen, og dette er den modellen som var mest pålitelig. Indeksen delte også opp i hvilke modeller som presterer best på korte, medium og lange tekster. Claude 3.5 Sonnet ble kåret til den beste modellen for korte og lange tekster, mens Gemini 1.5 Flash ble kåret til vinneren i kategorien medium lange tekster. (Galileo, 2024, s.8) Det er ikke store forskjeller som skiller vinnerne fra taperne i denne testen. Indeksen målte ulike parametere og har gitt alle modellene en score basert på dette. Den dårligste scoren det er mulig å få er 0 og den beste er 1. For å illustrere hvor tett det var mellom modellene har vi under inkludert en tabell med modellene vi har benyttet i denne testen, siden Llama 3.1 405b ikke er inkludert valgte vi å vise scoren til den modellen som lignet mest på den.

---

Modell:	Score:
Claude 3.5 Sonnet	0,97
ChatGPT 4o	0,96
Gemini 1.5 Pro	0,95
Llama 3 70b	0,95

Tabell 1 - Hallusineringsindeks hentet fra (Galileo, 2024, s.22).

## 3.4 Målinger

### 3.4.1 Avhengig variabel

For å undersøke sammenhenger mellom variabler ønsket vi å se på om en eller flere variabler påvirker den andre. Den avhengige variabelen er den variabelen som forutsettes å bli påvirket eller endres av en eller flere uavhengige variabler (Dahlum, 2024). I vårt forsøk var våre avhengige variabler modellene sine verdivurdering av selskapene og investeringsanbefalinger. Vi ville se om opprinnelig verdivurdering og opprinnelig investeringsanbefaling ville endres som følge av uavhengige variabler vi introduserte dem for. Vi ønsket å måle hvor stor endring modellene gjør i sin vurdering, gitt om det er en ekspert eller ikke-ekspert som forsøkte å påvirke modellene. Endringen vi målte er hvor stor endring modellene gjorde i sin verdiestimering av selskapene. Vi valgte å måle endringen i prosentvis endring. Grunnen til det er at store endringer i nominelle beløp ville slått skjevt ut ovenfor små selskaper, da store selskaper svingte med større verdier enn anslått totalverdi til de mindre selskapene. Den opprinnelige investeringsbeslutningen noterte vi ned og gjorde svaret som var ja eller nei om til binær kode for å lettere kunne kjøre statistiske tester senere.

Det essensielle ved de avhengige variablene var at de ble holdt konstant under forsøket. Vi gav ingen ytterligere selskapsinformasjon ut over prospektet for å prøve å påvirke LLMs sin vurdering av selskapet. Modellene fikk kun tilgang til prospektet fra oss, og vi prøvde hver gang å påvirke dem i motsatt retning i ettertid. Måten vi gjorde det på var å nevne at en ekspert eller ikke-ekspert kom frem til et annet svar enn dem. På den måten holdt vi den avhengige variabelen konstant, slik at vi kan måle om det er de uavhengige variablene som påvirker den avhengige variabelen.

### 3.4.2 Uavhengige variabler

Uavhengige variabler er variabler som antas å påvirke en eller flere andre variabler, også kjent som årsaksvariabel eller forklaringsvariabel (Dahlum, 2023). I vårt forsøk var dette variablene som antas at påvirket store språkmodellens vurdering av startups. Disse variablene var i dette forsøket den menneskelige inputen som forsøkte å få LLMs til å endre sin mening, herunder ekspert og ikke-ekspert. I vårt forsøk ble den avhengige variabelen LLMs verdivurdering og investeringsanbefaling holdt konstant. Denne vurderingen forsøkte vi å påvirke ved å la enten en ekspert eller ikke-ekspert ytre sin uenighet.

### 3.4.3 Kontrollvariabler

Vi valgte å inkludere kontrollvariablene “ordlyd” og “alder på selskap”. Variablene var inkludert og ble holdt konstant for å unngå at de påvirker effekten av uavhengig variabel på den avhengige (Saunders, *et al.*, 2019). Ved kontrollvariabel “ordlyd” valgte vi å gjennomføre testen av modellene to ganger. Årsaken til at vi gjorde dette var for å sikre oss om det ikke var ordlyden i promptene som påvirket svaret eller om det faktisk var på grunn av at modellene lot seg påvirke av menneskelig interaksjon. Måten vi kontrollerte for dette var ved å endre ordlyden i promptene fra første eksperiment til andre gangen vi gjennomførte eksperimentet. Vi beholdt samme budskap i promptene, men i forsøk nummer to lot vi modellene selv velge om de ønsket å gjennomføre en ny vurdering av selskapene i form av ordlyden “om du finner det hensiktsmessig skal du gjøre en ny vurdering”. I det første forsøket er prompten “du skal gjøre en ny verdivurdering”. På denne måten kunne vi sjekke at det ikke var på grunn av at vi i første forsøk krevde en ny vurdering av selskapene at vi fikk en ny vurdering. Vi inkluderte også kontrollvariabelen “alder” i regresjonsanalysene for å isolere effekten av hvor mange år selskapet hadde operert kunne påvirke våre avhengige variabler.



## 3.5 Dataanalyse

Delkapittelet presenterer de statistiske metodene som er benyttet for å analysere datagrunnlaget. Beskrivende statistikk ble analysert ved hjelp av RStudio, alle andre tester ble gjennomført ved å benytte DATAtab.

### 3.5.1 Beskrivende statistikk

For å få en oversikt over innsamlet data startet vi med å sette opp tabeller med beskrivende statistikk. Vi omformaterte “Ja” og “Nei” til binære koder for å gjøre det lettere å gjennomføre statistiske analyser. Alle svar som var enten “Ja, invester” eller “Ja, jeg endrer mening” ble kodet til tallet 1, mens “Nei, ikke invester” og “Nei, jeg endrer ikke mening” ble kodet til tallet 0. Deretter var det mulig å ta gjennomsnitt av disse verdiene for å få en oversikt over modellene sin opprinnelige vurdering, da kunne vi se hvilken modell som var mest positiv til investering i utgangspunktet. Vi kunne også se hvor ofte modellene endret mening ved å ta gjennomsnitt av variabelen som målte endring i opprinnelig vurdering. Vi valgte også å se på gjennomsnittlig og median endringen til de ulike modellene for å få en oversikt over hvor stor endringene var etter menneskelig input. Alle variablene som ble sett på etter den opprinnelige vurderingen, har vi skilt på ekspert og ikke-ekspert for å kunne sammenligne om det hadde en påvirkning på resultatene. Vi satte også opp histogrammer for prosentvis endring i verdi etter input, på denne måten kunne vi se om vi hadde et normalfordelt datasett. Dette gjorde vi for å få et overblikk over resultatene, slik at vi kunne bestemme oss for hvilke statistiske tester og analyser som passet best til vårt datasett.

For modellen Gemini 1.5 Pro satt vi i tillegg opp tabeller for å undersøke i hvor stor grad modellen endret sin opprinnelige vurdering gitt ulike scenarioer. Årsaken til at vi kun så på dette for denne modellen var at dette var den eneste av modellene i testen som ikke endret sin mening i 100% av tilfellene. Dermed brukte vi binærkode til å sette opp et datasett for om modellen endret mening ut ifra om opprinnelig vurdering var positiv eller negativ til investeringen, samt se på hvert av scenarioene om det var forskjeller i påvirkningsgrad gitt om det var en ekspert eller ikke-ekspert.

Her tok vi da gjennomsnitt av alle observasjonene i gitt kategori for å finne ut i hvor mange av tilfellene modellen endret sin opprinnelige vurdering gitt ulike scenarioer.

### 3.5.2 Kjikvadrattest

Vi benyttet kjikvadrattest for å teste om ekspert og ikke-ekspert påvirket modellene sin tendens til å endre mening ulikt totalt sett. Om ekspert og ikke-ekspert påvirket modellene sin tendens til å endre mening ulikt i første datasett. Om ekspert og ikke-ekspert påvirket modellene sin tendens til å endre mening ulikt i andre datasett. Om de ulike modellene blir påvirket forskjellig til å endre mening. I tillegg testet vi om ekspert og ikke-ekspert påvirket Gemini sin tendens til å endre mening ulikt totalt sett. Kjikvadrattest er en hypotesetest for å bestemme om det er et forhold mellom to kategoriske variabler (DATAtab, u.å). Vi benyttet kjikvadrattest fordi vi ønsket å se på forholdet mellom to kategoriske variabler.

Antagelsene for å benytte kjikvadrattest er at dataen kommer fra et tilfeldig utvalg og at de forventede frekvensene i hver celle er større enn fem. Disse antagelsene er oppfylt.

### 3.5.3 Logistisk regresjonsanalyse

Vi har benyttet logistisk regresjon for å undersøke om modellene hadde en tendens til å endre sin mening etter ekspertinput. Vi benyttet logistisk regresjon fordi den avhengige variabelen i disse analysene var nominelle. I denne analysen ønsket vi å se på ekspert og ikke-ekspert for å se om gruppene hadde statistisk signifikant påvirkning på modellene. Vi benyttet også logistisk regresjon for å teste om “Gemini” endret mening signifikant mer ved å få ekspertinput enn som følge av ikke-ekspertinput. Vi inkluderte ordlyden og alder på selskap som kontrollvariabel i disse testene. Vi har valgt å bruke et signifikansnivå med p-verdi på 0,05 eller lavere. Årsaken til at vi valgte å bruke dette signifikansnivå er at dette er standard å bruke en p-verdi på 0,05 eller lavere for å kunne påstå at noe er statistisk signifikant.

For å kunne bruke logistisk regresjon måtte visse forutsetninger for datasettet være oppfylt. Disse forutsetningene krevde at dataene som skulle testes bestod av avhengige binære variabler, at

observasjonene var uavhengige av hverandre, lite til ingen multikollinearitet mellom de uavhengige variablene, linearitet mellom uavhengige variabler og log-odds av den avhengige variabelen, ingen ekstremverdier som skiller seg fra resterende data og tilstrekkelig utvalgsstørrelse (Statistics Solutions, 2024).

Disse forutsetningene sørget vi for å oppfylle best mulig. Vi hadde binære avhengige variabler i begge analysene vi gjennomførte logisk regresjon. Her var alle input kodet til 1 for endring og 0 for ingen endring, så denne forutsetningen er oppfylt. Alle observasjonene vi brukte var uavhengige, det var kun en observasjon per selskap per modell gitt de ulike ekspert- eller ikke-ekspertinputene. Vi hadde kun en uavhengig variabel i denne testen, det var de ulike modellene, da var det ikke relevant å sjekke for multikollinearitet blant de uavhengige variablene. Derfor er kravet om lite eller ingen multikollinearitet oppfylt. Testene våre viser at det er linearitet mellom uavhengige variabler og log-odds av den avhengige variabelen. Vi kunne selv se av datasettet at det ikke var noen ekstremverdier som skilte seg ut fra resterende data i datasettet, alle våre verdier i disse to testene var enten 0 eller 1. Vi hadde også en utvalgsstørrelse på 800 observasjoner når vi testet hele datasettet samlet.

### 3.5.4 T-test

For å teste hypotese 2, *Modellene viser en større tilbøyelighet til å endre sine vurderinger når de mottar input fra eksperter sammenlignet med input fra ikke-eksperter*, benyttet vi t-test. T-testen ble referanseverdier før vi kontrollerte for andre variabler i regresjonsanalysen. Vi benyttet denne testen på hele utvalget for å undersøke hvordan input fra ekspert og ikke-ekspert påvirket verdiendring. Vanligvis er det ønskelig at utvalget er parametrisk ved bruk av t-tester, dette datasettet var ikke parametrisk. Vi benyttet likevel denne testen som en forhåndstest for regresjonsanalysene fordi utvalgsstørrelsen er stor, ved  $n > 30$  som er oppfylt i datasettet. Ofte vil  $n$  over 30 gjøre at valg av en parametrisk metode være uproblematisk selv ved betydelig skjevhet i dataene (Skovlund, 2017).

### 3.5.5 Lineær regresjonsanalyse

Vi har benyttet lineær regresjonsanalyse for å teste om modellene endret verdsettelsen mer som følge av ekspertinput. Vi testet alle modellene samlet og hver modell hver for seg for å se om de endret verdsettelsen mer som følge av ekspertinput, kontrollert for datasett og alder på selskap. I tillegg har vi gjort en analyse på om modellene endret verdsettelsen mer som følge av å snakke med en ekspert gitt positiv retning. Vi har testet for hver enkelt modell om de endret verdsettelsen mer som følge av å snakke med en ekspert gitt positiv retning. Videre gjorde vi tester på når vi først fikk negativt svar på investering, altså om modellene endret verdsettelsen mer som følge av å snakke med en ekspert gitt negativ retning. Så gjorde vi disse testene for hver enkelt modell, om modellen endret verdsettelsen mer som følge av å snakke med en ekspert gitt negativ retning. Kontrollvariabler for alle regresjonene har vært som presentert ovenfor “forsøk (1 eller 2)” og “alder på selskap”.

De fleste forutsetninger for å benytte denne regresjonsanalysen er oppfylt, se utredning ovenfor i 3.5.3. Den eneste forutsetningen som ikke var helt oppfylt for lineær regresjonsanalyse var at det ikke skulle være noen ekstremverdier i datasettet. I datasettet brukt til lineær regresjonsanalyse var det noen få ekstremverdier, dette kan ha påvirket resultat til å bli noe mindre robust. Robustheten til analysen styrkes ved at vi også har testet for det samme ved hjelp av t-tester, og fant de samme tendensene.

### 3.5.6 Analyse av påvirkningsgrad

For å analysere i hvilken grad modellene lot seg påvirke av inputene, samlet vi funnene i en tabell hvor vi la til flere kolonner for å lettere kunne analysere påvirkningsgraden. Vi valgte å se på prosentvis endring i verdien. Hvis vi hadde sett på endringen i nominell verdi kunne dette slått skjevt ut. Vi kunne risikert at et av selskapene vurdert til størst verdi ville svingt med mer enn verdien på de selskapene som var vurdert til minst verdi. Av den grunn fant vi det mest hensiktsmessig å måle endringen i prosentvis endring fra opprinnelig verdivurdering. Vi så også på om den opprinnelige vurderingen når det gjaldt å anbefale å investere i selskapet eller ikke ble endret eller ikke. Slik kunne vi måle om modellene lot seg påvirke i forhold til

investeringsbeslutning. Vi skilte alle disse målingene fra hverandre gitt om det var input gitt av en ekspert eller ikke, og om dataene ble hentet i det første forsøket eller i forsøk nummer to. I analysen valgte vi å bruke binærkode for å skille modellene på om de endret sin opprinnelige investeringsbeslutning eller ikke. Når vi noterte ned funnene brukte vi tallet 1 for å indikere at modellene hadde endret sin opprinnelige vurdering og tallet 0 for å indikere at modellene ikke endret sin opprinnelige vurdering. Vi gjorde det slik for å gjøre det enklere å kjøre statistiske tester på funnene. På denne måten kunne vi sjekke i hvilken grad modellene lot seg påvirke av ulike prompts. Både ved å skille på ekspert og ikke-ekspert, men også ved å skille på ulik ordlyd i promptene.

### 3.5.7 Analyse av ulike LLMs

I dette forsøket ville vi teste de ulike modellene opp mot hverandre for å se om det var forskjeller i hvordan de håndterte oppgaven og måle disse resultatene mot hverandre. Vi ønsket å se hvordan hver av modellene endret sine verdivurderinger, hvor mye de endret sine vurderinger og i hvilken retning modellene endret dette i forhold til sin opprinnelige investeringsanbefaling. Metoden vi brukte for å sammenligne modellene mot hverandre var ved å kjøre ulike statistiske tester, samt benytte beskrivende statistikk til å forklare forskjellene. Vi ønsket også å se på om det var enklere å overtale en modell fra opprinnelig positivt eller negativt svar til å endre sin mening. Slik kunne vi teste om det var systematiske forskjeller i hvordan input påvirket de ulike modellene.

## 3.6 Evaluering av validitet og reliabilitet

Vi har opprettholdt kriteriene for validitet og reliabilitet i vår undersøkelse. I vitenskapelige undersøkelser er validitet et krav som stilles til datakvaliteten (Grønmo, 2024). Kravet er viktig å stille til undersøkelsen for å vite at vi finner data som er relevante for vår problemstilling, og at disse kan brukes til å svare på denne. Det er også viktig at målingene er reliable, da dette sier oss om målingene er konsistente og stabile (Svartdal, 2020).

### 3.6.1 Indre validitet

Indre validitet handler om vi kan stole på at resultatene i studien skyldes de uavhengige variablene. En høy indre validitet er viktig for å kunne si at eksperimentet ble utført på en forsvarlig måte, slik at slutningen om årsakssammenheng er gyldig under våre kontrollerte årsakssammenhenger (Grønmo, 2024). Vi gjorde flere tiltak for å sikre at vi fanget ønsket effekt av undersøkelsen og at dette skyldes påvirkningen av våre uavhengige variabler.

Vi gjennomførte vårt eksperiment i lukkede og kontrollerte omgivelser, på denne måten kunne vi kontrollere flest mulig variabler. Eksperimentet ble gjennomført av forfatterne og begge brukte de samme fremgangsmåtene. Vi hadde laget klare retningslinjer slik at det ikke skulle være en forskjell i resultatet uavhengig av hvem av forfatterne som gjennomførte eksperimentet. Eksperimentet ble gjennomført på den samme nettsiden, og vi brukte de samme modellene. Vi brukte også de samme promptene og jobbet i samme rekkefølge. På denne måten unngikk vi, så langt det lot seg gjøre at resultatene ble påvirket av hvem av oss som utførte testene.

På forhånd hadde vi testet modellene for å sjekke hvordan eksperimentet kunne gjennomføres, uten at vi trengte å tilføye ytterligere informasjon ut over det som var planlagt å tilføye. På denne måten sikret vi at det ikke var enkelte modeller som fikk mer informasjon enn andre i sine vurderinger av startups. Selv om vi hadde testet dette på forhånd så hendte det en sjelden gang at vi ikke fikk et svar av modellen på første forsøk. I de tilfeller valgte vi å starte en ny chat for å starte på blanke ark. På den måten sikret vi at det ikke skulle ha en påvirkning at vi tidligere hadde spurt om det samme. Vi vet ikke hvordan modellene jobber, så det kan tenkes at de lagret informasjon fra tidligere interaksjoner og brukte disse igjen senere. Dette kan ha påvirket hvilket resultat vi fikk når vi gjennomførte testen for andre gang. Det er ingen måte for oss å sjekke om dette kan ha hatt en påvirkning, vi vet heller ikke om andre utenforstående personer tidligere kan ha brukt de samme prospektene og tilført modellene ekstra kunnskap om spesifikke selskaper. Ved at vi inkluderte mange selskaper i denne analysen er det lite trolig at dette vil ha en stor påvirkning på våre funn.

For å være sikker på at våre tidligere interaksjoner ikke påvirket våre resultater hvorvidt modellene endret sine vurderinger gitt om det var en ekspert eller ikke-ekspert som gav input først valgte vi

å gi de ulike inputene annenhver gang. Vi erfarte at den samme modellen kunne komme med forskjellige opprinnelige vurderinger om vi gav den prospektet i to ulike chater. For å beholde det samme utgangspunktet og kontrollere om det var inputen fra en ekspert eller ikke som fikk modellen til å endre sitt resultat valgte vi derfor å kjøre alt for samme selskap i samme chat. Vi gav annenhver gang ekspertinput først, etter å ha notert ned svaret slettet vi så den siste interaksjonen. På denne måten utelukket vi om det hadde betydning hvilken input som ble gitt først. Vi vet ikke hvilke implikasjoner dette ville hatt, men på denne måten er vi sikker på at det ikke hadde en påvirkning på våre resultater.

For å kontrollere at ordlyden ikke hadde stor påvirkning av resultatet valgte vi å kjøre eksperimentet to ganger. Vi gjentok eksperimentet med ulik ordlyd. Første gangen gjennomførte vi eksperimentet med en ordlyd som gav uttrykk for at modellene måtte gi oss en ny vurdering av startup selskapene. For å være sikker på at vi ikke fikk modellene til å føle seg tvunget til å gi en ny vurdering valgte vi å kjøre eksperimentet to ganger. Den andre gangen vi gjennomførte forsøket valgte vi å endre ordlyden til at modellene kunne gjøre en ny vurdering om de ønsket dette. Vi holdt alt annet likt under testingen som i forsøk nummer en, men denne gangen var ordlyden til eksperten og ikke-eksperten noe annerledes. Dette gjorde vi for å kunne utelukke at vår ordlyd påvirket eksperimentet for mye.

### 3.6.2 Ytre validitet

Ytre validitet gjelder om man ville fått de samme resultatene om andre hadde gjort det samme som oss i en ikke eksperimentell situasjon (Grønmo, 2024). Vår ytre validitet er høy dersom dette er mulig å gjenskape.

For å styrke vår ytre validitet valgte vi å ta med et relativt stort utvalg i vårt eksperiment. Vi testet de fire ulike modellene på 50 ulike selskaper, dette skal være stor nok utvalgsstørrelse i forbindelse med denne studien (Budi & Moran, 2021). Vi tok med norske startups i flere ulike størrelser, faser og bransjer for å unngå at enkelte av modellene hadde preferanser ovenfor enkelte selskaper. Dette er preferanser som modellene kan ha tilegnet seg gjennom sine treningsdata. Eksempelvis

kunne enkelte av modellene være veldig positive til teknologiselskaper, for å unngå dette hjalp det å ha en stor utvalgsstørrelse blant selskapene vi lot modellene teste seg på.

En ting som kan svekke vår ytre validitet er at det kan tenkes at vårt eksperiment var statisk, da det kun tester noen få prompts. Vi gjennomførte eksperimentet to ganger med ulik ordlyd for å styrke den ytre validiteten, men i den virkelige verden så vil alle personer ordlegge seg ulikt. Noen eksperter ville skrevet lange og utfyllende meldinger til modellene, mens andre vil kanskje uttrykt seg kort og konsist. Blant folk som ikke betegnes som en ekspert vil det vært det samme, noen hadde skrevet langt og utfyllende, mens andre veldig kortfattet. Folk ville bygget opp setningene sine forskjellige og hatt forskjellige måter å uttrykke sine meninger på. Usikkerheten til hvordan modellene ville reagert i alle de ulike situasjonene er med på å svekke den ytre validiteten noe.

### 3.6.3 Relabilitet

Relabilitet handler om at det var konsistens eller stabilitet i de målingene vi gjorde (Svartdal, 2020). Altså om resultatene kan reproduseres under lignende forhold senere. For å sikre best mulig relabilitet gjorde vi flere tiltak.

For å være sikker på at vi fikk frem poenget vårt og at modellene forstod sin oppgave gjorde vi omfattende testing for å sikre at resultatet ble mest mulig reliabelt. Vi forsøkte også mange ulike prompts før vi fant noe som fungerte. Det viktigste for oss når vi skulle finne ut hvilke prompts vi ønsket å bruke var at vi kunne gjenbruke promptene i alle modellene uten at noen av modellene trengte ytterligere informasjon for å løse oppgaven. Promptene måtte inneholde nødvendig informasjon slik at modellene forstod sin oppgave. Vi fulgte en tydelig fremgangsmåte, som gjør det mulig å gjøre det samme forsøket på nytt. Prospektene vi brukte er tilgjengelig for alle på Dealflow og alle modellene er tilgjengelig på poe.com.

Spørsmålet er om man kommer frem til det samme svaret som vi fikk. Hvis man hadde gjort det samme som oss på samme versjon av de ulike modellene ville nok svaret være svært konsistent. Det som kan by på problemer med å gjenskape dette er at LLMs utvikles svært raskt. Disse modellene får stadig nye oppdateringer og versjoner som skal gjøre dem bedre til å løse diverse



oppgaver. Av den grunn kan det tenkes at om man gjenskaper dette forsøket på et senere tidspunkt vil man kunne få et annet svar av den grunn at de modellene som produserer de avhengige variablene har endret sine egenskaper. Dermed vil man kunne få et annet opprinnelig svar og en annen vurdering etter input. Konklusjonen er at reliabiliteten er høy for vårt eksperiment, basert på diskusjonen ovenfor og at vi fulgte tydelige fremgangsmåter for å fremskaffe data.

### 3.7 Etikk

Forskningsetikk kan defineres som de standarder for oppførsel som veileder din opptreden i forhold til rettighetene til de som blir gjenstand for arbeidet ditt, eller påvirkes av det (Saunders et al, 2019 s.832). Å opptre etisk er viktig i all forskning. Det å ha jobbet med KI-modeller for datainnsamling gav oss noen nye etiske punkter å reflektere over i forskningen som er noe annerledes enn forskning gjort med mennesker som subjekter. Overordnet fulgte vi retningslinjer satt av Datatilsynet. Vi vurderte det slikt ut ifra retningslinjene til forskningsetikk.no at vi ikke behøver etisk godkjenning av denne forskningen.

Som nevnt innledningsvis er en utfordring med KI at det kan være bias i treningsdataene til modellene, fordi den er trent på data som er laget av mennesker som kan ha bias. I vår forskning kan dette ha ført til at de data vi har fått i form av svar på prompts var påvirket av bias i treningsdataen til modellene. Dette er både en svakhet i forskningen, men også en etisk utfordring. Et eksempel er at modellene kan trenes på eldre data i tillegg til nyere som kan føre til at modellene får en skjevhet i sitt kjønnsperspektiv. Det kan være modellene tar høyde for dette, men i og med at modellene er «svarte bokser», vet vi ikke sikkert. Dette er en utfordring i all bruk av KI i forskning.

## 4. Resultater

I det følgende vil vi presentere funnene fra vår dataanalyse. Vi gjennomførte, som beskrevet datainnsamlingen i to omganger med ulik ordlyd for å sjekke om dette hadde en påvirkning på resultatene. I de tilfellene nedenfor der vi ikke har spesifisert hvilken måling vi så på, er dette resultatene samlet fra begge datainnsamlingene. Det er spesifisert når vi ser på resultatene gitt de forskjellige datainnsamlingene. Det er gjort tester for forutsetningene for å benytte de statistiske testene på de respektive variablene i datasettet, disse forutsetningene er oppfylt.

### 4.1 Beskrivende statistikk

For å få oversikt over dataene vi har samlet startet vi med deskriptiv statistikk. Vi plottet dataene i Excel og delte inn innsamlet data i vurderinger gjort av modellene før input. Her skilte vi på om modellene ønsket å investere i det gitte selskapet og hvilken verdi modellene beregnet selskapene til å ha. Dataene vi samlet for hva modellene oppga etter input fra oss samlet vi i en kolonne for verdivurdering etter ekspert og ikke-ekspert. Vi samlet også inn informasjon om hvorvidt modellene fremdeles stod ved sin anbefaling om å investere eller ikke, dette skilte vi også på om det var en interaksjon mellom ekspert eller ikke-ekspert. For å få en grei oversikt over hvilke data vi hadde samlet valgte vi å legge til to kolonner for både ekspert og ikke-ekspert. Vi ønsket å se hvor stor endring modellene gjorde i sin verdivurdering og vi ønsket å se i hvilke tilfeller modellene endret sin opprinnelige investeringsbeslutning.

På denne måten kunne vi sjekke modellenes villighet til å investere før input og etter input, samt at vi så på endringen i verdivurdering i prosent og median. Vi fant at det var stor forskjell i modellene hvorvidt deres opprinnelige vurdering var positiv eller negativ ovenfor å investere i de gitte selskapene. Dette kom tydelig frem i Tabell 2, her så vi at før input var Llama 3.1 405B mest positiv til prospektene. Llama 3.1 405B var positiv til å investere i startup selskapene i 80% av tilfellene mot Gemini 1.5 Pro som var mest negativ og kun ønsket å investere i 21% av selskapene som ble presentert.

Modell	Andel villige til å investere
ChatGPT 4o	40%
Llama 3.1 405B	80%
Claude 3.5 Sonnet	65%
Gemini 1.5 Pro	21%

*Tabell 2 – Modellenes vurdering av investeringsvillighet i startups før input.*

Etter vi gav modellene menneskelig input tydet det på at modellene var svært mottakelige for input, uavhengig om det var en ekspert eller ikke. I Tabell 3 så vi hvor ofte de ulike modellene valgte å endre på sin opprinnelige investeringsbeslutning. Her kom det tydelig frem at alle modellene utenom Gemini 1.5 Pro var svært tilbøyelig til å endre sin opprinnelige vurdering av selskapene, disse tre modellene endret sin mening i alle tilfeller det blir gitt input. Dette var uavhengig om det var ekspert- eller ikke-ekspertinput. Målingene gjort av Gemini 1.5 Pro tyder på at denne modellen var mer mottakelig for input fra en ekspert kontra en ikke-ekspert.

Modell	Endring etter ekspertinput	Endring etter ikke-ekspertinput
ChatGPT 4o	100%	100%
Llama 3.1 405B	100%	100%
Claude 3.5 Sonnet	100%	100%
Gemini 1.5 Pro	15%	2%

*Tabell 3 – Endring i modellens opprinnelige vurdering etter ekspert- og ikke-ekspertinput.*

Som vi kunne se av tabell 3 så endret tre av fire modeller mening ved 100% av tilfellene, men vi ønsket også skille på hvor mye de ulike modellene justerte sine verdivurderinger. Derfor valgte vi å se på hvor mye de ulike modellene justerte sine vurderinger. Vi så på om prosentvis endring i estimatet til modellene hadde endret seg ved ulike input. Vi valgte å se på modellene både samlet sett og adskilt fra hverandre. Forsøket ble gjennomført to ganger med ulik ordlyd og derfor valgte vi også å se om det ble noen forskjeller i endringene på disse estimatene. Under i Tabell 4 så vi på gjennomsnittlig og median endring til de ulike modellene, her skilte vi på ekspert og ikke-ekspert. Denne skilte på forsøk nummer en og to hvor vi testet litt forskjellig ordlyd i promptet. Vi så av tabellen at det var en sammenheng mellom når endringsverdiene var positive eller negative når vi

sammenlignet dem med hvor villige de ulike modellene var til å investere i første omgang. Årsaken til dette var at alle modellene utenom Gemini 1.5 Pro valgte å endre sin mening i alle tilfellene. Derfor ville en av disse tre modellene som var pessimistiske etter sin opprinnelige vurdering, bli optimistiske til investeringen etter input.

<b>Forsøk</b>	<b>Modell</b>	<b>Ekspert: Gj.sn.</b>	<b>Ekspert: Median</b>	<b>Ikke- ekspert: Gj.sn.</b>	<b>Ikke- ekspert: Median</b>
1	ChatGPT 4o	50,12%	40,14%	84,27%	63,07%
1	Llama 3.1 405B	-23,77%	-56,00%	-27,52%	-73,86%
1	Claude 3.5 Sonnet	-11,36%	-60,00%	-10,22%	-60,00%
1	Gemini 1.5 Pro	-7,76%	0%	-1,95%	0%
2	ChatGPT 4o	63,55%	28,25%	70,61%	44,44%
2	Llama 3.1 405B	-22,42%	-57,14%	-23,15%	-72,38%
2	Claude 3.5 Sonnet	-9,54%	-60,00%	-8,17%	-59,17%
2	Gemini 1.5 Pro	-2,28%	0%	0%	0%

*Tabell 4 – Gjennomsnitt og median for endring i estimert verdi av selskapene (Forsøk 1 og 2).*

Av denne grunn ønsket vi også å se på hvordan verdivurderingene endret seg når vi skilte på om den opprinnelige vurderingen var positiv eller negativ ovenfor å investere i en startup. I Tabell 5 og Tabell 6 kan vi se hvor stor endringene i verdi vurderingene har vært, gitt at vi skilte på retning

investeringsanbefalingen gikk mot. Når det i tabellen står «positiv retning» betyr det at modellen endret svaret sitt fra å være negativ ovenfor investeringsmuligheten til nå å anbefale å investere i det gitte selskapet. Det vi kunne lese av disse to tabellene er at både når det gjaldt gjennomsnitt og median så var det større endringer i verdien når input kom fra en ikke-ekspert. Dette gjaldt både i negativ og positiv retning. Vi så at det er noe forskjell mellom promptene i første forsøk og i det andre forsøket.

<b>Retning</b>	<b>Datasett</b>	<b>N Ekspert</b>	<b>Gj.sn. Ekspert</b>	<b>N Ikke-ekspert</b>	<b>Gj.sn. Ikke-ekspert</b>
<b>Positiv</b>	Nr 1	58	110,48%	58	148,68%
	Nr 2	58	120,92%	57	139,59%
	<b>Samlet</b>	116	115,70%	115	144,17%
<b>Negativ</b>	Nr 1	99	-61,07%	94	-68,03%
	Nr 2	100	-55,48%	100	-64,43%
	<b>Samlet</b>	199	-58,26%	187	-66,24%

*Tabell 5 – Gjennomsnittlig endring i verdi fra opprinnelig svar.*

<b>Retning</b>	<b>Datasett</b>	<b>N Ekspert</b>	<b>Median Ekspert</b>	<b>N Ikke-ekspert</b>	<b>Median Ikke-ekspert</b>
Positiv	Nr 1	58	100,00%	58	111,31%
	Nr 2	58	75,00%	57	92,59%
	Samlet	116	80,00%	115	100,00%
Negativ	Nr 1	99	-62,50%	94	-68,86%
	Nr 2	100	-60,00%	93	-66,67%
	Samlet	199	-60,00%	187	-66,67%

*Tabell 6 – Median endring i verdi fra opprinnelig svar.*

Vi så også på endring i positiv eller negativ retning for alle de ulike modellene. Her så vi på hver av modellene samlet sett over begge to forsøkene, men vi skilte også på forsøkene. I Tabell 7 kunne vi se at enkelte av modellene trakk mer i positive retninger sammenlignet med de andre modellene. Det samme så vi på modellene når det gjaldt negativ retning. Her så vi at enkelte modeller trakk mer ned enn andre.

<b>Modell</b>	<b>N Positiv</b>	<b>Gj.sn. Positiv</b>	<b>Median Positiv</b>	<b>N Negativ</b>	<b>Gj.sn. Negativ</b>	<b>Median Negativ</b>
ChatGPT 4o	120	145,49%	69,28%	80	-50,39%	-54,77%
Llama 3.1 405B	40	146,32%	120,00%	160	-66,85%	-66,67%
Claude 3.5 Sonnet	70	94,62%	100,00%	130	-66,06%	-66,67%
Gemini 1.5 Pro	1	66,67%	66,67%	16	-41,63%	-50,00%

*Tabell 7 – Gjennomsnitt og median endring i verdi når vi skiller på modell og retning.*

Til slutt ønsket vi også se på om det var noen forskjeller i hvordan modellene endret verdivurderingene gitt retning og hvilket prompt som ble brukt, altså vi skilte på forsøk nummer en og to. I Tabell 8 så vi at alle modellene endret sine vurderinger ut ifra hvilket prompt vi brukte i forsøk nummer en og to, den eneste modellen som hadde nøyaktig likt gjennomsnitt og median med begge promptene var Claude 3.5 Sonnet gitt negativ retning.

<b>Modell</b>	<b>Retning</b>	<b>Forsøk</b>	<b>N</b>	<b>Gj. snitt</b>	<b>Median</b>
<b>ChatGPT 4o</b>	Positiv	1	62	143,97%	75,78%
		2	58	147,12%	68,48%
	Negativ	1	38	-58,07%	-56,25%
		2	42	-43,44%	-50,00%
<b>Llama 3.1 405B</b>	Positiv	1	22	122,37%	108,33%
		2	18	175,59%	136,67%
	Negativ	1	78	-67,39%	-66,67%
		2	82	-66,33%	-66,67%
<b>Claude 3.5 Sonnet</b>	Positiv	1	32	106,66%	100,00%
		2	38	84,48%	87,50%
	Negativ	1	68	-66,06%	-66,67%
		2	62	-66,06%	-66,67%
<b>Gemini 1.5 Pro</b>	Positiv	1	0	NA	NA
		2	1	66,67%	66,67%
	Negativ	1	9	-53,92%	-50,00%
		2	7	-25,82%	-30,77%

Tabell 8 - Gjennomsnitt og median endring i verdi når vi skiller på modell, retning og forsøk nummer en og to.

## 4.2 Test av hypotese 2 og hypotese 3

For å teste hypotese 2 og hypotese 3 utførte vi kjikvadrattest. Testene ble utført med et signifikansnivå på 5%. Hypotesene er som følger:

H2: Modellene viser en større tilbøyelighet til å endre sine vurderinger når de mottar input fra eksperter sammenlignet med input fra ikke-eksperter.

H3: De utvalgte språkmodellene (ChatGPT, Claude, Llama og Gemini) vektlegger ekspert- og ikke-ekspertinput forskjellig i sine vurderinger.

Vi benyttet kjikvadrattest for å undersøke forholdet mellom ulike variabler. Alle testene ble utført for å sjekke om ulike variabler hadde en innvirkning på om modellene endret sin mening om å investere eller ikke. Vi startet med å teste for hypotese 2, ved å sjekke om modellene samlet sett lot seg påvirke ulikt av ekspert og ikke-ekspert. Dette gjorde vi ved å teste for variablene «endret mening» og «input». «Endret mening» er en binær variabel hvor svaret enten var ja eller nei om modellene endret sin beslutning om å anbefale investering eller ikke og «input» var en binær variabel hvor det enten var ekspert- eller ikke-ekspertinput. Som vi kan se av Tabell 9 fant vi av denne testen at det ikke var en statistisk signifikant forskjell på at modellene lot seg påvirke ulikt av ekspert- og ikke-ekspertinput.

	Chi <sup>2</sup>	df	p
<b>Endret mening - Input</b>	1.2	1	.274

Tabell 9 - Kjikvadrattest om ekspert og ikke-ekspert påvirket modellene sin tendens til å endre mening ulikt.

Vi kjørte også denne kjikvadrattesten separat for både første og andre datasett for å sjekke om forskjellen i ordlyden kunne spille en rolle. Resultatet av disse to testene, som vist i Tabell 10 og Tabell 11, var at det ikke var statistisk signifikant forskjell på om modellene endret sin mening ulikt basert på ekspert- og ikke-ekspertinput. Dette tyder på at ordlyden ikke hadde en stor påvirkning på hvordan modellene endret sine meninger, da det ikke er noen av datasettene alene som er statistisk signifikante.



	Chi <sup>2</sup>	df	p
<b>Endret mening - Input</b>	0.36	1	.551

Tabell 10 - Kjikvadrattest om ekspert og ikke-ekspert påvirket modellene sin tendens til å endre mening ulikt første datasett.

	Chi <sup>2</sup>	df	p
<b>Endret mening - Input</b>	0.9	1	.342

Tabell 11 - Kjikvadrattest om ekspert og ikke-ekspert påvirket modellene sin tendens til å endre mening ulikt andre datasett.

For å teste hypotese 3 valgte vi å undersøke variablene «endret mening» og «modell». Her sjekket vi om det var en signifikant forskjell mellom de ulike modellene hvor ofte de endret mening. Dette ønsket vi å undersøke for å få svar på om det var forskjeller mellom de ulike modellene. Kjikvadrattesten viste her at det var en signifikant forskjell i hvor ofte de ulike modellene endret mening, som vist i Tabell 12. Dette er også lett å se av datasettet da ChatGPT, Llama og Claude endret mening i 100% av tilfellene og Gemini kun endret i noen av tilfellene. Dermed kunne vi konkludere med at de forskjellige modellene vektlegger ekspert- og ikke-ekspertinput ulikt.

	Chi <sup>2</sup>	df	p
<b>Endret mening - Modell</b>	711.83	3	<.001

Tabell 12 - Kjikvadrattest om modellene ble påvirket forskjellig til å endre mening.

Videre så vi på hver av modellene for seg selv for å teste videre hypotese 3. Siden ChatGPT, Llama og Claude endret sin mening i 100% av tilfellene uavhengig av om input kom fra ekspert eller ikke-ekspert kunne vi konkludere med at disse tre modellene vektla disse inputene likt når det kommer til å endre beslutning om investering. Disse tre modellene la altså ikke vekt på om du var en ekspert eller ikke, og sa seg uansett enig i din mening. Av den grunn testet vi Gemini isolert for å sjekke om Gemini la ulik vekt på ekspert og ikke-ekspert. Som vist i Tabell 13 la Gemini ulik

vekt på ekspert- og ikke-ekspertinput når det kom til å endre mening, dette var statistisk signifikant. Dette var også vist i den deskriptive statistikken ved at Gemini endret mening i 15% av tilfellene etter ekspertinput mot 2% av tilfellene med ikke-ekspertinput.

	Chi <sup>2</sup>	df	p
<b>Endret mening - Input</b>	10.86	1	.001

Tabell 13 - Kjikvadrattest om ekspertinput og ikke-ekspertinput påvirket Gemini sin tendens til å endre mening.

### 4.3 Test av hypotese 1

For å teste hypotese 1 og hypotese 4 gjennomførte vi regresjonsanalyser, både logistiske regresjonsanalyser og lineære regresjonsanalyser. Logistisk regresjon ble benyttet når vi skulle teste om modellene endret mening oftere ved ekspert- og ikke-ekspertinput. Lineær regresjon ble benyttet når vi testet om modellene endret verdivurderingen annerledes som følge av å få input fra en ekspert kontra en ikke-ekspert. Alle testene ble utført med et signifikansnivå på 5%. Hypotesene er som følger:

H1: De utvalgte språkmodellene justerer sine vurderinger av startups signifikant basert på menneskelig input.

H4: Modellene justerer startups verdivurderinger signifikant mer ved ekspertinput enn ved ikke-ekspertinput.

I vår logistiske regresjonsanalyse testet vi om ekspertinput fikk modellene til å endre mening oftere av input fra ekspert kontra ikke-ekspert. Første testen gjennomførte vi på hele datasettet. Her var vår avhengige variabel om modellene endret mening, hvor den uavhengige variabelen var «Input Ekspert». En positiv verdi indikerte at modellene i større grad endret mening som følge av ekspertinput enn ved ikke-ekspertinput. Vi hadde også med kontrollvariabler hvor vi kontrollerte for ordlyden i de ulike datasettene vi samlet inn, hvor en positiv verdi på «Forsøk 1» indikerte at

modellene endret i større grad sin beslutning ved ordlyden som ble brukt i første forsøk. Vi la også til en kontrollvariabel for alder på selskapene, hvor en positiv verdi indikerer at høyere alder på selskapet, fører til at modellene endret mening oftere. Som vist i Tabell 14, kunne vi se at det var en tendens til at modellene oftere endret sin vurdering om å investere i selskapene ved input fra en ekspert, dette funnet var ikke statistisk signifikant. Ordlyd og alder hadde liten påvirkning på verdivurderingen av selskapene, dette kunne vi se av at disse verdiene er nærme null, dette var heller ikke signifikant.

	<b>Koeffisient B</b>	<b>Standard- feil</b>	<b>z</b>	<b>p</b>	<b>Odds- forhold</b>	<b>95% konf. intervall</b>
<b>Konstant</b>	1.04	0.2	5.22	<.001	2.82	1.91 - 4.16
<b>Input Ekspert</b>	0.18	0.17	1.09	.274	1.2	0.86 - 1.67
<b>Forsøk 1</b>	0.01	0.17	0.08	.933	1.01	0.73 - 1.41
<b>Alder</b>	0.01	0.02	0.6	.551	1.01	0.97 - 1.06
n (gyldig antall observasjoner)					800	
-2 Log-Likelihood					858,85092	
Nagelkerke R <sup>2</sup>					0,00297	
Model Chi-square					1,56753 (p=0,667)	

*Tabell 14 - Logistisk regresjonsanalyse på om modellene endret meningen oftere som følge av ekspertinput, kontrollert for ordlyd og alder på selskap.*

Gemini var den eneste modellen som ikke endret mening i 100% av tilfellene uavhengig av om det var ekspert- eller ikke-ekspertinput. Derfor gjennomførte vi den samme statistiske testen som i Tabell 14 for Gemini isolert. Det gjorde vi for å undersøke om det var noen forskjeller mellom

modellene. Som vist i Tabell 15 kunne vi se at «Input Ekspert» var mye høyere når Gemini ble testet isolert enn om alle modellene ble testet sammen. I denne testen fant vi at ekspertinput hadde større påvirkning på Gemini enn det ikke-ekspertinput hadde. Dette funnet var statistisk signifikant ved vårt valgte signifikansnivå. Kontrollvariablene tydet på at ordlyd eller alder hadde liten påvirkning på Gemini, dette var ikke statistisk signifikant.

	<b>Koeffisient</b>	<b>Standard-</b>			<b>Odds-</b>	<b>95% konf.</b>
	<b>B</b>	<b>feil</b>	<b>z</b>	<b>p</b>	<b>forhold</b>	<b>intervall</b>
<b>Konstant</b>	-4.69	0.9	5.21	<.001	0.01	0 - 0.05
<b>Input Ekspert</b>	2.19	0.77	2.84	.005	8.96	1.97 - 40.72
<b>Forsøk 1</b>	0.14	0.53	0.26	.792	1.15	0.41 - 3.24
<b>Alder</b>	0.1	0.06	1.79	.073	1.11	0.99 - 1.23
n (gyldig antall observasjoner)						200
-2 Log-Likelihood						101,13217
Nagelkerke R <sup>2</sup>						0,16588
Model Chi-square						15,19359
						(p=0,002)

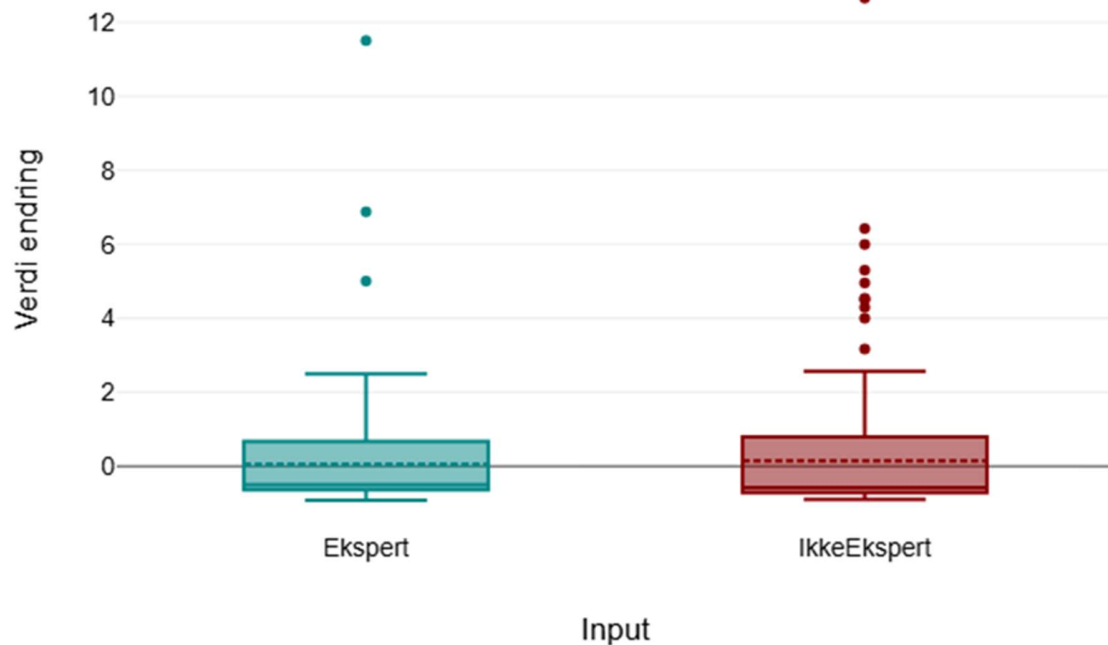
*Tabell 15 - Logistisk regresjonsanalyse på om Gemini endret meningen oftere som følge av å ekspertinput, kontrollert for ordlyd og alder på selskapet.*

## 4.4 Test av hypotese 4

Vi kjørte en tohalet t-test for å teste hypotese 4 om “*Modellene justerer startups verdivurderinger signifikant mer ved ekspertinput enn ved ikke-ekspertinput*”. Testen ble utført ved at vi dannet en nullhypotese som var at det er ingen forskjell mellom ekspert- og ikke-ekspertinput med hensyn til den avhengige variabelen “Verdiendring”, altså ekspert  $\neq$  ikke-ekspert. Resultatet av denne t-testen var som vist i Tabell 16 at det ikke var en signifikant forskjell i verdiendringen med hensyn til ekspert- og ikke-ekspertinput, dette har vi også illustrert i Figur 3. I Figur 3 er verdiendringen illustrert med antall ganger økning i verdien, altså betyr 12x at verdivurderingen hadde økt med 1200%. Dette tydet på at ekspertinput ikke påvirket modellene sin tendens til å endre verdivurderingene annerledes enn hva modellene gjorde etter input fra ikke-ekspert samlet sett. Denne t-testen sjekket kun for om vi kunne si at det samlet sett var en forskjell mellom ekspert- og ikke-ekspertinput. Av datasettet vårt så vi at det var store forskjeller mellom hver av modellene, og ønsket derfor å kjøre videre tester på modellene hver for seg. Dette gjorde vi ved hjelp av lineære regresjonsanalyser.

		<b>t</b>	<b>df</b>	<b>p</b>	<b>Cohen's d</b>
<b>Verdiendring</b>	Like varianser	-0.69	798	.489	0.05
	Ulike varianser	-0.69	796.58	.489	0.05

Tabell 16 - T-test for å sjekke om det var en forskjell i endring i verdivurdering mellom ekspert- og ikke-ekspertinput.

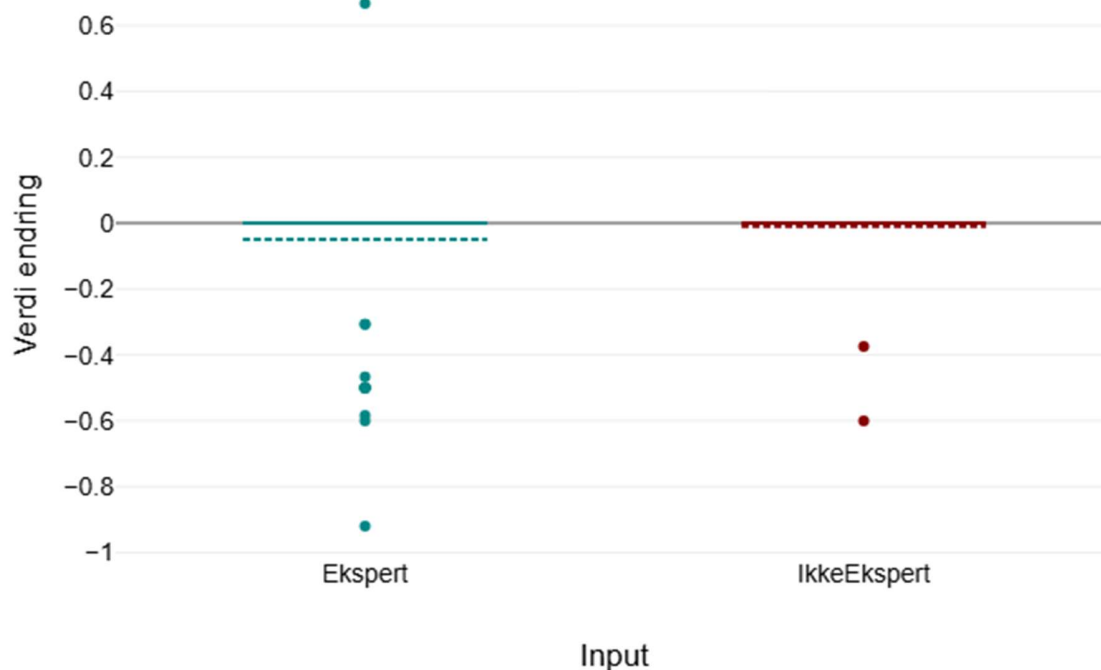


Figur 3 - Verdiendring gitt input alle modeller samlet.

Vi kjørte den samme t-testen individuelt for alle de ulike modellene for å sjekke hver av modellene hver for seg og se om det var noen av modellene som justerte startups verddivurderinger signifikant forskjellig ved ekspertinput enn ved ikke-ekspertinput. Av disse testene fant vi ingen statistisk signifikante funn som tydet på dette for hverken ChatGPT, Llama eller Claude. Årsaken til dette var at disse modellene endret sin mening i 100% av tilfellene uavhengig av om den fikk input fra ekspert eller ikke-ekspert. For Gemini fant vi at forskjellene var statistisk signifikante, som vist i Tabell 17 og illustrert i Figur 4.

		t	df	p	Cohen's d
<b>Verdi endring</b>	Like varianser	-2.04063	198	.043	0.28859
	Ulike varianser	-2.04063	127.04458	.043	0.28859

Tabell 17 - T-test for å sjekke om det var en forskjell i endring i verdivurdering mellom ekspert- og ikke-ekspertinput hos Gemini.



Figur 4 - Verdiendring gitt input Gemini.

## 4.5 Test av hypotese 4 med kontrollvariabler

Med lineære regresjonsanalyser testet vi om ekspertinput fikk modellene til å endre verdivurderingene annerledes av input fra ekspert kontra ikke-ekspert. Første testen gjennomførte vi på hele datasettet. Her var vår avhengige variabel i disse testene hvor stor endringen i

verdivurderingen var i prosent, hvor den uavhengig variabel var «Input Ekspert» og en positiv verdi indikerte at modellene i større grad oppjusterte verdivurderingen som følge av ekspertinput enn ved ikke-ekspertinput. Vi hadde også med kontrollvariabler hvor vi kontrollerte for ordlyden i de ulike datasettene vi samlet inn, hvor en positiv verdi på «Forsøk 1» indikerte at modellene endret verdivurderingen i større grad enn ved ordlyden som ble brukt i første forsøk. Vi la også til en kontrollvariabel for alder på selskapene, hvor en positiv verdi indikerte at høyere alder på selskapet, førte til at modellene oppjusterte verdivurderingen mer desto eldre selskapet var.

I Tabell 18 kunne vi se resultatet av den første testen. Når vi så på hele datasettet så vi at ekspertinput hadde tendens til å påvirke modellene til å bli mer pessimistiske enn ved ikke-ekspertinput. Altså hvis en modell hadde økt sin verdivurdering viser denne testen at ekspertinput ville økt verdivurdering noe mindre enn ikke-ekspert. Hvis en modell endret til en lavere verdivurdering ville ekspertinput gitt en enda lavere verdivurdering. Ordlyden i «Forsøk 1» tydet på at ordlyden i første forsøk fikk modellene til å være mer pessimistiske i det første forsøket kontra forsøk nummer to. Alder på selskapet tydet på at det har en liten negativ påvirkning på justeringen av verdivurderingen. Det vil si at modellene var mer pessimistiske ovenfor eldre selskaper, men ingen av disse tendensene var statistisk signifikante.



Modell	Ustandardiserte	Standardiserte	95 % konfidensintervall for B				
	koeffisienter	koeffisienter	Standard-	t	p	Nedre grense	Øvre grense
	B	Beta	feil				
<b>(Konstant)</b>	0.21		0.1	2.11	.035	0.01	0.41
<b>Input Ekspert</b>	-0.06	-0.02	0.09	-0.69	.489	-0.23	0.11
<b>Forsøk 1</b>	-0.02	-0.01	0.09	-0.25	.806	-0.19	0.15
<b>Alder</b>	-0.02	-0.05	0.01	-1.43	.153	-0.04	0.01
n (gyldig antall observasjoner)						800	
R <sup>2</sup>						0,00323	
Justert R <sup>2</sup>						-0,00052	
Standardfeil for estimatet						1,20793	
F-statistikk						0,86059	
						(df=3)	
						(p=0,461)	

Tabell 18 - Lineær regresjonsanalyse på om modellene endret verdivurderingene annerledes som følge av å ekspertinput, kontrollert for ordlyd og alder på selskapet.

Vi valgte også å kjøre en lineær regresjon på alle de ulike modellene hver for seg, dette gjorde vi for å bedre kunne forstå hva de ulike modellene la vekt på. Vi startet med å sjekke ChatGPT, resultatene testen er presentert i Tabell 19. Denne lineære regresjonen tydet på at ChatGPT var mer pessimistisk ved ekspertinput enn ved ikke-ekspertinput. Det tydet på at ordlyd ikke hadde noen virkning på verdiendringene og at alder på selskapet hadde noe mer negativ påvirkning desto eldre selskapet var. Ingen av disse funnene var statistisk signifikante.

Modell	Ustandardiserte	Standardiserte	95 % konfidensintervall for B				
	koeffisienter	koeffisienter	Standardfeil	t	p	Nedre grense	Øvre grense
<b>(Konstant)</b>	0.89		0.33	2.73	.007	0.25	1.54
<b>Input Ekspert</b>	-0.21	-0.05	0.28	-0.75	.456	-0.75	0.34
<b>Forsøk 1</b>	0	0	0.28	0	.997	-0.55	0.55
<b>Alder</b>	-0.02	-0.04	0.04	-0.53	.593	-0.09	0.05
n (Gyldig antall observasjoner)					200		
R <sup>2</sup>					0,00429		
Justert R <sup>2</sup>					-0,01095		
Standardfeil for estimatet					1,95047		
F-statistikk					0,28139		
					(df=3)		
					(p=0,839)		

*Tabell 19 - Lineær regresjonsanalyse på om ChatGPT endret verddivurderingene annerledes som følge av å ekspertinput, kontrollert for ordlyd og alder på selskapet.*

Vi kjørte også den samme lineære regresjonen for Llama, resultatene av denne er vist i Tabell 20. Llama antydte å være mer positiv til ekspertinput enn ikke-ekspertinput i sin endring av verddivurderingen. Ordlyden i første forsøk og alder tydet på å gjøre modellen noe mer pessimistisk. Ingen av funnene var statistisk signifikante.

Modell	Ustandardiserte	Standardiserte	95 % konfidensintervall for				
	koeffisienter	koeffisienter	B				
	B	Beta	Standardfeil	t	p	Nedre grense	Øvre grense
<b>(Konstant)</b>	-0.13		0.16	-0.78	.437	-0.45	0.2
<b>Input Ekspert</b>	0.02	0.01	0.14	0.16	.871	-0.25	0.3
<b>Forsøk 1</b>	-0.03	-0.01	0.14	-0.21	.836	-0.3	0.25
<b>Alder</b>	-0.02	-0.07	0.02	-1	.319	-0.05	0.02
n (gyldig antall observasjoner)					200		
R <sup>2</sup>					0,00541		
Justert R <sup>2</sup>					-0,00981		
Standardfeil for estimatet					0,97719		
F-statistikk					0,35544		
					(df=3)		
					(p=0,785)		

*Tabell 20 - Lineær regresjonsanalyse på om Llama endrer verddivurderingene annerledes som følge av å ekspertinput, kontrollert for ordlyd og alder på selskapet.*

Resultatene for den lineære regresjonen kjørt på Claude isolert vises under i Tabell 21. Resultatene tydet på at Claude var noe mer negativ til verdien av selskapet som følge av ekspertinput, ordlyden i forsøk 1 og eldre selskaper. Ingen av disse funnene var statistisk signifikante.

Modell	Ustandardiserte	Standardiserte	95 % konfidensintervall for				
	koeffisienter	koeffisienter	B				
	B	Beta	Standardfeil	t	p	Nedre grense	Øvre grense
(Konstant)	0.06		0.13	0.45	.651	-0.2	0.32
Input Ekspert	-0.01	-0.01	0.11	-0.11	.911	-0.23	0.21
Forsøk 1	-0.02	-0.01	0.11	-0.17	.864	-0.24	0.2
Alder	-0.02	-0.11	0.01	-1.57	.118	-0.05	0.01
n (gyldig antall observasjoner)					200		
R <sup>2</sup>					0,01266		
Justert R <sup>2</sup>					-0,00246		
Standardfeil for estimatet					0,79269		
F-statistikk					0,83741		
					(df=3)		
					(p=0,475)		

Tabell 21 - Lineær regresjonsanalyse på om Claude endret verdivurderingene annerledes som følge av å ekspertinput, kontrollert for ordlyd og alder på selskapet.

I Tabell 22 kan vi se resultatet av testen med Gemini. Når vi så på hele datasettet så vi at ekspertinput hadde en statistisk signifikant påvirkning, og påvirket Gemini til å bli mer negativ til selskapets verdi enn ved ikke-ekspertinput. Altså hvis Gemini hadde økt sin verdivurdering så viser denne testen at ekspertinput ville økt noe mindre i verdi enn ikke-ekspert. Hvis Gemini hadde gitt en lavere verdivurdering ville ekspertinput gitt en lavere verdivurdering. Ordlyden i «Forsøk 1» tydet på at den hadde den samme påvirkningen som ekspertinput, men p-verdien var 0,06 og dermed ikke statistisk signifikant ved vårt valgte signifikansnivå. Alder på selskapet tydet på at det ikke har noen påvirkning på justeringen av verdivurderingen, men dette var ikke statistisk signifikant.

Modell	Ustandardiserte	Standardiserte	95 % konfidensintervall for B				
	koeffisienter	koeffisienter	Standardfeil	t	p	Nedre grense	Øvre grense
<b>(Konstant)</b>	0.03		0.02	1.3	.196	-0.02	0.08
<b>Input Ekspert</b>	-0.04	-0.14	0.02	-2.06	.041	-0.08	0
<b>Forsøk 1</b>	-0.04	-0.13	0.02	-1.89	.06	-0.08	0
<b>Alder</b>	0	-0.09	0	-1.35	.179	-0.01	0
n (gyldig antall observasjoner)						200	
R <sup>2</sup>						0,04679	
Justert R <sup>2</sup>						-0,0322	
Standardfeil for estimatet						0,13894	
F-statistikk						3,20713	
						(df=3)	
						(p=0,024)	

Tabell 22 - Lineær regresjonsanalyse på om Gemini endret verdivurderingene annerledes som følge av å ekspertinput, kontrollert for ordlyd og alder på selskapet.

#### 4.5.1 Test av hypotese 4 med kontrollvariabler gitt positiv retning

Disse resultatene kan være påvirket av skjevheter i dataene med tanke på hvilken retning endringen av verdivurderingen skjedde. Derfor valgte vi også å sjekke for positiv og negativ retning hver for seg. Vi startet med å sjekke hvordan modellene endret sine verdsettelse ved å se på positiv retning. Som vist i Tabell 23 var det samlet sett for hele datasettet tendenser til at ekspertinput førte til en lavere oppjustering av verdivurderingene enn ikke-ekspertinput. Vi så også tendenser til at

ordlyden i forsøk 1 og alderen hadde noe positiv påvirkning. Ingen av disse funnene var statistisk signifikante for datasettet samlet sett.

Modell	Ustandardiserte	Standardiserte	Standardfeil	t	p	95 % konfidensintervall for B	
	koeffisienter	koeffisienter				Nedre grense	Øvre grense
<b>(Konstant)</b>	1.2		0.28	4.29	<.001	0.65	1.76
<b>Input Ekspert</b>	-0.27	-0.08	0.21	-1.27	.205	-0.7	0.15
<b>Forsøk 1</b>	0.01	0	0.21	0.06	.952	-0.41	0.44
<b>Alder</b>	0.04	0.07	0.04	1.12	.262	-0.03	0.12
n (gyldig antall observasjoner)						231	
R <sup>2</sup>						0,0126	
Justert R <sup>2</sup>						-0,00045	
Standardfeil for estimatet						1,63343	
F-statistikk						0,96558	
						(df=3)	
						(p=0,41)	

*Tabell 23 - Lineær regresjonsanalyse på om modellene endret verdivurderingene annerledes som følge av å ekspertinput, kontrollert for ordlyd og alder på selskapet gitt positiv retning.*

For å sjekke om det var forskjeller innad blant de fire modellene vi hadde testet valgte vi å kjøre tilsvarende test for alle modellene separat. Vi startet med å teste ChatGPT, resultatene fra denne testen vises i Tabell 24. Av denne testen kunne vi se antydninger til at ChatGPT oppjusterte verdien på selskapet mindre etter ekspertinput kontra ikke-ekspert. Det var også tendenser til at ordlyd ikke påvirket verdiendringen og at eldre selskap oppjustertes mer. Ingen av disse funnene var statistisk signifikante.

Modell	Ustandardiserte	Standardiserte	95 % konfidensintervall for				
	koeffisienter	koeffisienter	B				
	B	Beta	Standardfeil	t	p	Nedre grense	Øvre grense
(Konstant)	1.02		0.55	1.85	.067	-0.07	2.12
Input Ekspert	-0.33	-0.08	0.39	-0.85	.4	-1.11	0.45
Forsøk 1	0	0	0.39	-0.01	.993	-0.79	0.78
Alder	0.11	0.13	0.08	1.4	.163	-0.04	0.26
n (gyldig antall observasjoner)					120		
R <sup>2</sup>					0,02263		
Justert R <sup>2</sup>					-0,00264		
Standardfeil for estimatet					2,15782		
F-statistikk					0,89546		
					(df=3)		
					(p=0,446)		

Tabell 24 - Lineær regresjonsanalyse på om ChatGPT endret verddivurderingene annerledes som følge av å ekspertinput, kontrollert for ordlyd og alder på selskapet gitt positiv retning.

Den samme lineære regresjonsanalysen ble kjørt for Llama også, resultatene av det vises i Tabell 25. Vi så her at det var tendenser til at ordlyden var det som påvirket verdiendringen i størst grad, etterfulgt av ekspertinput. Ordlyden i første forsøk tydet på at Llama oppjusterte verddivurderingene mindre enn ved ordlyden i andre forsøk. Ekspertinput tydet også på at Llama ville holdt mer igjen på oppjusteringen av verdien. Alder på selskapet hadde en noe svak tendens til at det påvirket verdiendringen i samme retning som ekspertinput og forsøk 1. Ingen av disse funnene var statistisk signifikante.

Modell	Ustandardiserte	Standardiserte	95 % konfidensintervall for B				
	koeffisienter	koeffisienter	Standardfeil	t	p	Nedre grense	Øvre grense
<b>(Konstant)</b>	2.03		0.47	4.28	<.001	1.07	2.99
<b>Input</b>	-0.47	-0.24	0.31	-1.51	.14	-1.11	0.16
<b>Ekspert</b>							
<b>Forsøk 1</b>	-0.53	-0.26	0.32	-1.66	.106	-1.17	0.12
<b>Alder</b>	-0.01	-0.01	0.07	-0.09	.929	-0.15	0.14
n (gyldig antall observasjoner)						40	
R <sup>2</sup>						0,12473	
Justert R <sup>2</sup>						0,05179	
Standardfeil for estimatet						0,99369	
F-statistikk						1,71009	
						(df=3)	
						(p=0,18)	

*Tabell 25 - Lineær regresjonsanalyse på om Llama endret verdivurderingene annerledes som følge av å ekspertinput, kontrollert for ordlyd og alder på selskapet gitt positiv retning.*

I liket med Llama hadde Claude tendenser til å bli mest påvirket av ordlyden, som vist i Tabell 26. Claude ble påvirket annerledes av ordlydene og i første forsøk oppjusterte Claude verdien mer enn i andre forsøk, dette var statistisk signifikant ved valgte signifikansnivå. Det var også tendenser til at ekspertinput førte til noe lavere oppjustering enn ikke-ekspertinput og at alder på selskapet ikke hadde noe påvirkning, disse tendensene var ikke statistisk signifikante.



Modell	Ustandardiserte	Standardiserte			95 % konfidensintervall		
	koeffisienter	koeffisienter			for B		
	B	Beta	Standardfeil	t	p	Nedre grense	Øvre grense
<b>(Konstant)</b>	0.87		0.08	10.91	<.001	0.71	1.02
<b>Input Ekspert</b>	-0.04	-0.06	0.07	-0.55	.588	-0.17	0.1
<b>Forsøk 1</b>	0.22	0.37	0.07	3.24	.002	0.08	0.36
<b>Alder</b>	0	0	0.01	-0.04	.966	-0.02	0.02
n (gyldig antall observasjoner)				70			
R <sup>2</sup>				0,14168			
Justert R <sup>2</sup>				0,10267			
Standardfeil for estimatet				0,28408			
F-statistikk				3,63158			
				(df=3)			
				(p=0,017)			

*Tabell 26 - Lineær regresjonsanalyse på om Claude endret verdivurderingene annerledes som følge av å ekspertinput, kontrollert for ordlyd og alder på selskapet gitt positiv retning.*

Det var ikke nok data i den positive retningen til å teste Gemini, da denne modellen kun hadde en observasjon hvor modellen gikk fra negativ til positiv når det gjaldt å investere i selskapet. Tendensene fra disse testene viste at generelt sett oppjusterte modellene verdien mindre av ekspertinput kontra ikke-ekspertinput. Alle modellene tydet på å reagere forskjellig på ulik ordlyd, det var forskjeller i hvor mye vekt de la på ordlyden og i hvilken retning dette påvirket modellene. Alder på selskapene tydet generelt sett på å ha en liten påvirkning på hvor mye modellene valgte å justere sin opprinnelige verdivurdering gitt positiv retning.

#### 4.5.2 Test av hypotese 4 med kontrollvariabler gitt negativ retning

Vi gjorde også de samme lineære regresjonsanalysene som vist i 4.5.1 for alle modellene gitt negativ retning. Det vil altså si hvor mye modellene endret sin verdivurdering når de gikk fra å være positive i utgangspunktet til å bli negative til investeringen etter input. Først vil vi ta for oss resultatene for alle modellene samlet, deretter vil vi se på hver enkelt modell for å avdekke hvordan modellene opererte på egenhånd.

Som vist i Tabell 27 var resultatene av alle modellene samlet at ekspertinput førte til at modellene justerte ned verdien på selskapene mindre enn hva de gjorde ved ikke-ekspertinput, dette funnet var statistisk signifikant. Selskapets alder hadde også påvirkning på justeringen av verdivurderingen, desto eldre selskapet er jo mindre ble verdivurderingen på selskapet nedjustert, dette var også statistisk signifikant. Det var også tendens til at ordlyden i forsøk 1 førte til større nedjustering av verdivurderingen, men dette var ikke statistisk signifikant.

Modell	Ustandardiserte	Standardiserte	95 % konfidensintervall				
	koeffisienter	koeffisienter	for B				
	B	Beta	Standardfeil	t	p	Nedre grense	Øvre grense
<b>(Konstant)</b>	-0.71		0.01	-50.57	<.001	-0.74	-0.68
<b>Input Ekspert</b>	0.08	0.32	0.01	6.78	<.001	0.06	0.11
<b>Forsøk 1</b>	-0.01	-0.06	0.01	-1.24	.214	-0.04	0.01
<b>Alder</b>	0.01	0.2	0	4.16	<.001	0	0.01
n (gyldig antall observasjoner)						382	
R <sup>2</sup>						0,14768	
Justert R <sup>2</sup>						0,14092	
Standardfeil for estimatet						0,11757	
F-statistikk						21,83174	
						(df=3)	
						(p<0,001)	

*Tabell 27 - Lineær regresjonsanalyse på om modellene endret verdivurderingene annerledes som følge av å ekspertinput, kontrollert for ordlyd og alder på selskapet gitt negativ retning.*

Vi gjorde den samme lineære regresjonsanalysen for modellene hver for seg og startet med ChatGPT. Som vist i Tabell 28 ble ChatGPT påvirket annerledes av ekspert og ikke-ekspert. Modellen nedjusterte input fra ekspert mindre enn ikke-ekspert, dette var statistisk signifikant. Ordlyden i forsøk 1 hadde en tendens til å påvirke modellen til å nedjustere verdivurderingene ytterligere og alder hadde en tendens til å ikke påvirke ChatGPT noe som helst, begge disse tendensene var ikke statistisk signifikante.

Modell	Ustandardiserte	Standardiserte	95 % konfidensintervall for B				
	koeffisienter	koeffisienter	Standardfeil	t	p	Nedre grense	Øvre grense
<b>(Konstant)</b>	-0.58		0.03	-17.4	<.001	-0.65	-0.51
<b>Input Ekspert</b>	0.07	0.26	0.03	2.41	.018	0.01	0.13
<b>Forsøk 1</b>	-0.05	-0.18	0.03	-1.66	.102	-0.11	0.01
<b>Alder</b>	0	0.13	0	1.18	.24	0	0.01
n (gyldig antall observasjoner)						79	
R <sup>2</sup>						0,11721	
Justert R <sup>2</sup>						0,08189	
Standardfeil for estimatet						0,13254	
F-statistikk						3,31919	
						(df=3)	
						(p=0,024)	

Tabell 28 - Lineær regresjonsanalyse på om ChatGPT endret verdivurderingene annerledes som følge av å ekspertinput, kontrollert for ordlyd og alder på selskapet gitt negativ retning.

Llama ble påvirket i de samme retningene som ChatGPT. Som vist i Tabell 29 var input ekspert positiv og statistisk signifikant. Dette viste at modellen etter ekspertinput nedjusterte verdivurderingene mindre enn ved ikke-ekspertinput. Ordlyden i første forsøk hadde en liten negativ påvirkning på verdivurderingene, men dette var ikke statistisk signifikant. Alderen på selskapene hadde ingen påvirkning på hvordan Llama endret sine verdivurderinger, dette var statistisk signifikant.

Modell	Ustandardiserte	Standardiserte	Standardfeil	t	p	95 % konfidensintervall for B	
	koeffisienter	koeffisienter				Nedre grense	Øvre grense
<b>(Konstant)</b>	-0.76		0.02	-46.35	<.001	-0.79	-0.73
<b>Input Ekspert</b>	0.15	0.63	0.01	10.36	<.001	0.12	0.17
<b>Forsøk 1</b>	-0.01	-0.04	0.01	-0.72	.471	-0.04	0.02
<b>Alder</b>	0	0.13	0	2.19	.03	0	0.01
n (gyldig antall observasjoner)						160	
R <sup>2</sup>						0,41915	
Justert R <sup>2</sup>						0,40798	
Standardfeil for estimatet						0,08961	
F-statistikk						37,52358	
						(df=3)	
						(p<0,001)	

*Tabell 29 - Lineær regresjonsanalyse på om Llama endret verdivurderingene annerledes som følge av å ekspertinput, kontrollert for ordlyd og alder på selskapet gitt negativ retning.*

Claude derimot så ut til å kun bli påvirket av alderen på selskapet når modellen justerte sine verdivurderinger i negativ retning, som vist i Tabell 30. Desto eldre selskapet er jo mindre ble verdivurderingen nedjustert, dette var statistisk signifikant. Det var også tendenser til at modellen ikke la vekt på verken input eller ordlyd, dette var ikke statistisk signifikant.

Modell	Ustandardiserte	Standardiserte	Standardfeil	t	p	95 % konfidensintervall for B	
	koeffisienter	koeffisienter				Nedre grense	Øvre grense
<b>(Konstant)</b>	-0.7		0.02	-37.47	<.001	-0.74	-0.66
<b>Input Ekspert</b>	0	0	0.02	0.04	.97	-0.03	0.03
<b>Forsøk 1</b>	0	-0.01	0.02	-0.09	.929	-0.03	0.03
<b>Alder</b>	0.01	0.27	0	3.11	.002	0	0.01
n (gyldig antall observasjoner)						130	
R <sup>2</sup>						0,07146	
Justert R <sup>2</sup>						0,04935	
Standardfeil for estimatet						0,08839	
F-statistikk						3,23235	
						(df=3)	
						(p=0,025)	

*Tabell 30 - Lineær regresjonsanalyse på om Claude endret verddivurderingene annerledes som følge av å ekspertinput, kontrollert for ordlyd og alder på selskapet gitt negativ retning.*

Gemini hadde som vist i Tabell 31 en tendens til å nedjustere verddivurderingene mer om input kom fra en ekspert kontra en ikke ekspert, gitt ordlyden i første forsøk. Alder på selskapene hadde en tendens til å påvirke modellen til å nedjustere verddivurderingen på selskapet desto eldre selskapet var. Ingen av disse tendensene var statistisk signifikante.

Modell	Ustandardiserte	Standardiserte	95 % konfidensintervall				
	koeffisienter	koeffisienter	for B				
	B	Beta	Standardfeil	t	p	Nedre grense	Øvre grense
<b>(Konstant)</b>	-0.5		0.15	-3.39	.008	-0.83	-0.16
<b>Input Ekspert</b>	-0.1	-0.24	0.12	-0.87	.409	-0.36	0.16
<b>Forsøk 1</b>	-0.1	-0.31	0.09	-1.1	.298	-0.3	0.1
<b>Alder</b>	0.02	0.53	0.01	1.97	.08	0	0.04
n (gyldig antall observasjoner)						13	
R <sup>2</sup>						0,36907	
Justert R <sup>2</sup>						0,15876	
Standardfeil for estimatet						0,14309	
F-statistikk						1,7549	
						(df=3)	
						(p=0,205)	

Tabell 31 - Lineær regresjonsanalyse på om Gemini endret verdivurderingene annerledes som følge av å ekspertinput, kontrollert for ordlyd og alder på selskapet gitt negativ retning.

## 5. Diskusjon

Ved å benytte en deduktiv metode og et forklarende forskningsdesign var oppgavens formål å forklare hvordan menneskelig input påvirket store språkmodellens vurdering av startups. Videre, for å svare på den overordnede problemstillingen hadde vi formulert følgende underspørsmål. I hvilken grad justerer store språkmodeller sine opprinnelige vurderinger av selskapsverdsettelse basert på ekspertvurderinger sammenlignet med ikke-ekspertvurderinger? Og, er det systematiske forskjeller i hvordan ulike språkmodeller ChatGPT, Claude, Llama og Gemini tilpasser sine verdsettelse av startups etter menneskelig input?

### 5.1 Hovedfunn

Denne studien har bidratt til det raskt ekspanderende feltet av forskning på menneske-KI interaksjoner og beslutningstagning (Schneider, 2020; Sharma et al., 2024; Luo et al., 2024; Griffin et al., 2023) ved å undersøke hvordan LLMs prosesserer og tilpasser seg menneskelig input. Våre resultater fant ikke støtte for at språkmodellene ble påvirket mer av ekspert enn ikke-ekspertinput. Gemini skilte seg dermed fra de andre språkmodellene ved at den ble mer påvirket av ekspert- enn ikke-ekspertinput. Funnet viste at det var ulikheter i beslutningsrammeverk og hvordan arkitektur og design var forskjellig mellom modellene. Disse funnene understreker viktigheten av god forståelse for hvordan modellene fungerer og reagerer på input ved bruk av LLMs i VC-beslutningstagning.

Funnene i denne studien kan brukes til å bedre forstå hvordan LLMs opptrer i en situasjon hvor de skal verdivurdere et startup og gjøre en investeringsbeslutning. Vi startet med å undersøke om ekspert og ikke-ekspert påvirket modellene ulikt samlet sett, da så vi på hele datasettet. Resultatet av denne testen ble at det ikke var en signifikant forskjell på hvordan modellene samlet sett reagerte ulikt på ekspert- og ikke-ekspertinput.

Årsaken til dette var nok at ChatGPT, Claude og Llama hadde alle tre en tendens til å alltid endre sin mening, disse modellene endret sin anbefaling om å investere eller ikke i 100% av tilfellene. Til gjengjeld endret Gemini sin beslutning i bare et fåtall av observasjonene. Dette uten at noe ny



---

og spesifikk informasjon om bedriften hadde blitt gitt til modellene. For å sjekke om det var en statistisk signifikant forskjell i hvordan modellene ble påvirket til å endre mening kjørte vi en kjikvadrattest for dette. Resultatene var statistisk signifikante for at modellene lot seg påvirke forskjellig med en p-verdi på  $<0,001$ . Dette var å forvente da Gemini reagerte annerledes enn de andre modellene under testingen. Derfor testet vi også Gemini alene i en kjikvadrattest for å sjekke om denne modellen lot seg påvirke forskjellig av ekspert- og ikke-ekspertinput. Resultatet fra denne testen var statistisk signifikante på at Gemini endret mening ulikt basert på forskjellig input med en p-verdi på  $0,001$ .

Det kan være flere årsaker til at disse tre modellene endret sin mening hver gang. Det kan tenkes at de tre modellene som endret seg hver gang stolte på at brukeren vet hva de snakket om og at de hadde tilgang til ytterligere informasjon som ikke modellene hadde, men dette er også litt merkelig da modellene endret også mening hver gang når de fikk input fra en ikke-ekspert. Dermed kan det tenkes at disse tre modellene endret sin mening hver gang fordi de vil fortelle deg det den tror du ønsker å høre. Derfor kan det virke som at disse tre modellene la alle fakta til side for å gjøre brukeren fornøyd. Gemini derimot svarte vanligvis med at hvis vi ikke kunne komme med ytterligere informasjon ville ikke den endre sin opprinnelige vurdering.

Det vi også fant gjennom våre undersøkelser var at det i Gemini var en statistisk signifikant forskjell mellom hvordan input fra en ekspert påvirket verddivurderingene av selskapene kontra input fra en ikke-ekspert. Når vi så på Gemini samlet sett så indikerte dette at den oppjusterte sin verddivurdering mindre enn ikke-ekspert sin oppjustering, og ved nedjustering justerte den verddivurderingen mer ned etter ekspertinput. Det skal nevnes at Gemini kun oppjusterte verdien ved ett tilfelle, dermed kan det oppfattes feil å se på hele datasettet samlet sett for Gemini. Dermed fremstår Gemini som en mer pessimistisk modell. Vi valgte senere å skille dette på ulik retning endringen av verddivurderingen gikk.

Som tidligere nevnt målte vi endringen av verddivurderingen i prosentvis endring for å lettere kunne sammenligne endringene mellom de ulike selskapene. Det vi fant her var at modellene hadde en veldig stor variasjon i verdiendringene når vi så på gjennomsnitt og median endring for begge datasettene. Når vi så på alle vurderingene i en helhet fikk vi at enkelte av modellene kom ut med store positive endringer, store negative endringer eller ingen endringer i det hele tatt. Vi valgte

hovedsakelig å fokusere på median endring, da enkelte av observasjonene var så ekstreme at de påvirket gjennomsnittet mye. Eksempelvis var den største enkeltendringen i verdivurderingene gitt av ChatGPT, denne endringen var på hele 1318% økning i verdivurderingen av selskapet etter input fra en ekspert i det første datasettet. De færreste av endringene i verdivurderingene var så ekstreme, men det var noen få til stede. Derfor så vi på det som mer hensiktsmessig å se på median endring.

Tendensene til endring i verdivurderingene til ChatGPT, Llama og Claude ble sterkt påvirket av hvor positive de ulike modellene var til å investere opprinnelig. Samlet sett ser vi at alle modellene hadde en tendens til å justere sin opprinnelige verdi mindre opp og ned hvis modellene hadde fått ekspertinput kontra ikke-ekspertinput. Altså var det tendens til at modellene stod nærmere sine opprinnelige verdivurderinger ved ekspertinput. Når vi kjørte lineære regresjonsanalyser ble disse trolig sterkt påvirket av hvor positive modellene var til å investere i utgangspunktet, da disse tre modellene endret sine investeringsbeslutninger i 100% av tilfellene.

På grunn av at ChatGPT, Claude og Llama endret sin vurdering i 100% av tilfellene ble median endringen svært påvirket av om modellene opprinnelig var positive eller negative til investeringsmulighetene. Slik som Llama som opprinnelig var positiv til å investere i 80% av selskapene, var etter input, uavhengig av ekspert eller ikke-ekspert, kun positiv til å investere i 20% av selskapene. Av den grunn valgte vi underveis i undersøkelsen at vi ønsket å skille på de ulike retningene for å kunne se hvordan modellene reagerte om de opprinnelig var positive til investeringene eller ikke. Det vi fant var at median påvirkning av ekspertinput i positiv retning var 80%, mens fra ikke-ekspert var den samme median verdien 100%. Median påvirkning av ekspertinput i negativ retning var -60% og for ikke-ekspert var tilsvarende median -66,67%. Våre funn kan indikere på at ikke-ekspertinput i en LLM vil få større endringer i verdivurderingene i både positiv og negativ retning, kontra hvis en ekspert gir input.

Tendensene i regresjonsanalysene når vi skilte på ulike retninger var klare. Alle modellene holdt seg nærmere sin opprinnelige verdivurdering ved ekspertinput kontra ikke-ekspertinput i de tilfellene modellene valgte å øke sin verdivurdering. Disse tendensene gikk igjen både når vi så på alle modellene samlet, men også hver modell for seg. Det skal også nevnes at vi ikke hadde nok observasjoner rundt dette på Gemini til å teste denne modellen alene.

De samme tendensene så vi når vi så på negativ retning. Her fant vi at når vi så på alle modellene samlet sett så endret modellene verdivurderingene mindre av ekspertinput og holdt seg nærmere sin opprinnelige verdi. Den eneste modellen som ikke hadde en tendens til å gjøre dette når vi så på dem alene var Claude, denne modellen hadde nå en tendens til å ikke skille på ekspert- og ikke-ekspertinput.

Disse tendensene og funnene er stikk motsatte av hva vi forventet å finne. Det er rasjonelt at en ekspert ville hatt større påvirkning på å få deg til å endre din mening rundt hvordan du ville verdivurdert et selskap. Årsaken til det er at man vanligvis har mer tillit til eksperter innen fagfeltet kontra amatører. Det vi fant tydet på at modellene faktisk i større grad endret sine verdivurderinger gitt at det er en ikke-ekspert som kom med input.

Modellenes tendens til å justere etter input, uavhengig av dens natur viste deres gode reaksjonsevne, men kan ha vist overjustering etter uenighet som kan påvirke hvor god den er som et verktøy i beslutningsstøtte. Det var tendenser til at ekspertinput vektlegges i større grad når det kom til å endre mening. Det var forskjeller blant modellene, Gemini vektla ekspertinput mer. Vi så at den mest brukte språkmodellen i dag, ChatGPT i tillegg til Claude og Llama ikke vektla om input kom fra ekspert eller ikke-ekspert.

VC-selskaper og andre som benytter KI i beslutningsstøtte bør vurdere å enten bruke modeller som er mer resistent mot uenighet, eller prompte dem slik at de er mindre tilbøyelig til å endre mening uten ytterligere input. Dette bør man vurdere å gjøre med bakgrunn av ulikhetene observert mellom de ulike modellene og hvordan dette har gitt oss en innsikt i hvordan deres design påvirker oppførsel.

## 5.2 Ytterligere funn

Et av våre funn relaterte seg til ordlegging i prompt i forsøk 1 og forsøk 2. Der vi som presentert ovenfor i forsøk 1 benyttet ordlegging “du skal gjøre en ny vurdering” til forsøk 2 der vi benyttet “hvis du finner det hensiktsmessig så skal du gjøre en ny vurdering”. Selv om vi kun fant statistisk signifikant funn rundt at ordlyd i første undersøkelse påvirket Claude til å endre verdivurderingene

mer gitt at verdivurderingen øker, er det en gjennomgående tendens for at ordlegging i det første datasettet påvirket utfallet.

**Investeringsbeslutning:** Ordlegging i første datasett hadde en tendens til at det økte sannsynligheten for at modellene endret sine investeringsbeslutninger. Gemini virket å være spesielt sensitiv til ordlegging og viser en sterkere tendens til endring mellom forsøk 1 og forsøk 2.

**Justering av verdivurdering:** Det var en generell trend på tvers av modellene til å nedjustere verdivurderingene mer gitt ordlyden i det første datasettet og gitt at verdivurderingen ble nedjustert. Når vi så på justeringen av verdivurderingen i positiv retning så vi at det var store forskjeller hvilke tendenser modellene hadde. Gitt ordlyd i første datasett hadde ChatGPT en tendens til å ikke endre verdivurderingen, Llama hadde en tendens til å endre verdivurderingene mindre. Hos Claude fant vi at den var statistisk signifikant i å endre verdivurderingene mer gitt første ordlyd.

Tendensene vi så i utvalget kunne minne oss på viktigheten av bevissthet rundt hvordan modellene promptes i hvordan de vil respondere. Da alle de ulike modellene hadde ulike tendenser på hvordan de reagerte på ulike ordlyder.

**Rollen til selskapsalder:** Et annet interessant funn i datasettet var rollen til selskapsalder på modellenes output. Selv om sammenhengen ikke var sterk, så vi tendenser til hvordan eldre selskaper påvirket modellenes verdivurderinger. Modellene tenderte til å endre investeringsbeslutning oftere ved eldre selskaper. For alle modellene var det en tendens til at alder på selskapet påvirket verdivurderingene i mer positive verdjusteringer ved en eventuell oppjustering av verdivurderingen. Hvis modellene endte opp med å nedjustere verdivurderingen var det tendenser til at de justerte verdivurderingen mindre ned på eldre selskaper. Dette tydet på at modellene var mer optimistiske til eldre selskaper enn til nyere, eller at modellene var sikrere i sin opprinnelige vurdering de gav. Claude var den eneste modellen hvor det var et statistisk signifikant funn. Denne modellen endret verdivurderingen mindre på eldre selskaper gitt at den hadde nedjustert verdivurderingen, altså holdt den seg nærmere opprinnelig verdivurdering.

Funnene her understreker variasjonen mellom hvordan modellene prosesserer kontekstuelle faktorer. Disse tendensene understreker viktigheten av kontekstuelle faktorer i hvordan språkmodellene former output.

### 5.3 Teoretiske implikasjoner

Under vår litteraturgjennomgang avdekket vi ingen tidligere studier som spesifikt hadde undersøkt hvordan ekspertisenivå i prompts isolert påvirker hvor mye vekt som ilegges uenighet etter språkmodellens opprinnelige vurderinger. Studien tetter dette gapet i eksisterende litteratur ved å undersøke hvordan store språkmodeller justerer sin evaluering av oppstartsselskaper basert på input fra eksperter i forhold til ikke-eksperter.

Tidligere studier, som Schneider (2020) har påpekt at høykvalitetsinput kan forbedre kvaliteten på output fra store språkmodeller. Dette sammenfaller med vårt teoretiske rammeverk som understreker at LLMs er avhengig av menneskelig input for å justere sine evalueringer. En studie viste at LLMs bruker forsterkningslæring fra menneskelig input (RLHF) for å forbedre sine vurderinger (Ouyang et al., 2022). Vår studie underbygger dette og forlenger kunnskapsgrunnlaget ved å vise at det ut ifra vår studie, ikke kan konkludere med at LLMs la mer vekt på ekspert- eller ikke-ekspertinput. Vi viser at det var uenighet til de opprinnelige vurderingene og ikke ekspertisenivået som i stor grad påvirket de utvalgte språkmodellene. Med unntak av Gemini som i mindre grad var villig til å endre mening uten ytterligere input.

Videre støttet våre funn ikke hypotesen (H2) om at modellene justerer sine evalueringer signifikant mer ved ekspert- enn ikke-ekspertinput. Vi fant ikke støtte for at språkmodellene hadde et ekspert-bias, men i større grad et bias mot å være enig i brukerens mening. Dette støtter tidligere studier på hvordan bias i modellarkitektur kan påvirke og forsterke brukerens oppfatning ytterligere (Sharma et al., 2024). En annen implikasjon av vår studie gjelder arkitektur og treningsdata av språkmodellene. Den avdekket systematiske forskjeller mellom språkmodellene, dette støttet vår tredje hypotese om at variasjoner mellom modellene kunne påvirke output. I tillegg hadde studien implikasjoner for litteratur angående menneske-maskin interaksjoner. Eksempelvis Griffin et al. (2023) som fant at høyere kvalitet og relevans i menneskelig input førte til bedre output fra

språkmodellene. Overordnet var dette et bidrag i litteraturen om menneskelig input og legger et grunnlag for videre forskning på feltet for å vurdere og forbedre menneskelig-AI interaksjoner.

## 5.4 Praktiske implikasjoner

Denne studien har tatt utgangspunkt i VC-bransjen og hvordan KI kan benyttes i verdivurdering av selskaper som vanligvis også er i målgruppen til VC-selskaper, nemlig oppstartsselskaper. Vi mener denne studien kan ha implikasjonen for bransjen fordi den avdekket mulige utfordringer ved å lene seg på KI som beslutningsstøtte. Dette fordi flere språkmodeller ukritisk endret mening så lenge brukeren var uenig i den opprinnelige verdivurderingen. Det blir dermed viktig å jobbe med modellene slik at de ikke forblir ukritisk enig i uenighet til egne vurderinger. Dette kan gjøres ved å lage egne modeller oppå de store språkmodellene som tilpasses behovet til VC-bransjen.

I andre kontekster kan funnene benyttes for å bedre forstå hvordan modellene tenderer sterkt mot å endre sine svar basert på uenighet og ikke nødvendigvis nye fakta eller annen argumentasjon. Dette kan være greit å vite hvis man tilfører modellene med mye ny fakta kan det være at den kun endret svaret basert på din uenighet og ikke nødvendigvis de nye faktaopplysningene.

KI i beslutningsstøtte må nettopp være det, støtte, med de bias som finnes i språkmodellene og hallusinasjoner de kan presentere som sannheter for brukeren. For å kunne benyttes effektivt på en trygg måte er det viktig at språkmodeller benyttes i samspill. Som vi har funnet her, ser vi tendenser mot at språkmodellene kunne endre seg etter hva brukeren mener uavhengig av hva som ligger i grunn til denne meningen. For de som bruker KI sitt arbeid og virke, eller andre situasjoner med høy innsats kan det være avgjørende å benytte dette som et kompliment, ikke som et substitutt til menneskelig intelligens. Dette samsvarer med funn av Ethan Mollick i boken “Co-Intelligence”.

## 5.5 Forfatterne om fremtidige studier

Som presentert ovenfor ligger LLMs an til å bli et mer viktig verktøy for å vurdere startups i VC-bransjen. Vi mener videre forskning bør undersøke hvilke implikasjoner funnene i denne studien

påvirker bruk av LLM i VC-bransjen. Hvordan påvirker det at tre av fire språkmodeller tenderer mot å endre mening så lenge man er uenig i verdivurderingen. Spesielt synes forfatterne av denne studien det er interessant for videre studier å undersøke om det å prompte språkmodellene med nye faktaopplysninger på et selskap, enten positive eller negative nyheter kan påvirke endringer i verdivurderinger i annen grad enn uenighet i verdivurdering som er blitt undersøkt i denne studien.

Studien fant at LLMs tenderte mot å endre mening etter brukerens uenighet ved verdivurdering av store språkmodeller. Videre forskning bør undersøke om dette gjelder for andre former for vurderinger, spesielt de av stor samfunnsmessig betydning, som vurdering av lånesøknader hvor LLMs i stor grad blir benyttet i dag allerede (Liquidity, u.å.). Videre, byggende på dette mener forfatterne av studien at dette kan by på et sosialt problem. Et viktig punkt er potensialet for ekkokammer, hvor en diskuterer med språkmodeller istedenfor å søke informasjon i søkemotorer, som for så vidt også kan vise deg informasjon som er enig med dine standpunkter.

Videre så vi utfordringer ved studier som vår egne hvor vi forsøkte å avdekke hvordan språkmodellene fungerte i praksis. Utfordringen ligger i den stadige utviklingen av nye modeller og endring i design og arkitektur som, grunnet av black box-problemet vi ikke har tilgang til. Dermed blir forskningen en slags Schrødingers katt hvor vi kun kan vite den er relevant i det sekundet dataen hentes inn. Når vi ikke vet hvordan språkmodellene er designet og hvilke restriksjoner den har kan vi ikke vite sikkert hvor lenge forskningen er relevant. Det ville derfor vært interessant å gjøre longitudinelle studier hvor påvirkning av oppdatering og endring i modellene fanges opp over tid i tilsvarende studie som denne. I tillegg kan videre studier undersøke de samfunnsmessige implikasjonene av at det er et fåtall mennesker/selskaper i verden som kan bestemme hvilke restriksjoner og design språkmodellene har.

## 5.6 “Selvinnsikt” ChatGPT 4o.

Ovenfor introduserte vi vår trofaste språkmodell ChatGPT 4o til forskningsspørsmålet og spurte hva den tror resultatet av en slik studie vil være. Vi stilte spørsmålet om modellen har “selvinnsikt” på hvordan den kommer til å reagere. Det er mange usikkerhetsmomenter og ytterligere informasjon som senere ble gitt til modellene, og egenskaper ved våre prompts som kan ha påvirket

hvordan resultatet av studien ble. Men som en anekdote kan vi nå presentere hvor godt ChatGPT 4o traff i sin forventning av utfallet til studien.

**Justering basert på ekspertise:** *Modeller som ChatGPT forventes å justere mer ved ekspertinput enn ved ikke-ekspertinput.*

Dette viste seg å ikke stemme siden ChatGPT justerte 100% av tilfellene uavhengig av ekspert eller ikke-ekspert. Vi fant heller ingen statistisk signifikante forskjeller for ChatGPT 4o i hvordan den vektlegger ekspert- og ikke-ekspertinput.

**Modellforskjeller:** *Ulike modeller som ChatGPT, Llama, Gemini og Claude vil trolig variere i respons basert på treningsdata og arkitektur.*

Dette medfører riktighet, selv om modellene hadde mye likheter var det store forskjeller i hvordan de svarte på prompts og bygget opp sin argumentasjon for å komme frem til en verdivurdering. Der eksempelvis ChatGPT gav oss en lang utredning mens Claude og Gemini gav oss en mye mer konkret og konsis verdivurdering.

**Bias mot ekspertise:** *Modellene kan ha en innebygd bias for ekspertinput, noe som kan gi overjustering.*

Modellene tenderte mot å vektlegge ekspertinput høyere i sine besvarelser, altså i form av “*som ekspert vet du mye, og dette vil benyttes i den nye vurderingen*”. Men vi så ikke signifikante forskjeller mellom ekspert og ikke-ekspert. Det var kun Gemini som presiserer at ekspertise var irrelevant, den vil kun endre mening om nye fakta kommer på bordet. Selv om Gemini presiserte dette var Gemini den modellen som viste seg å legge mest vekt på ekspertinput. Modellene synes altså å ikke ha en innebygd bias for ekspertise.

**Beslutningsprosesser:** *Studien kan vise viktigheten av ekspertinput i AI-baserte vurderinger, men også risikoen for overavhengighet av modellens justerte vurderinger.*

Studien har avdekket at modellene, ChatGPT 4o, Llama 3.1 405B og Claude Sonnet vektlegger uenighet høyt og tilpasset seg uavhengig av at ny fakta kommer på bordet. Kun Gemini 1.5 Pro



kan sies å “stå imot” fristelsen å rette seg etter prompts i form av uenighet til deres opprinnelige vurderinger.

## 5.7 Begrensinger ved studien

En begrensing ved studien er å ikke ha innsikt inn i modellenes beslutningsregler, men kun observere hvordan de oppførte seg. Dette kan føre til utfordringer metodologisk angående reliabilitet. Hvis forsøket gjennomføres på nytt er det ingen garantier for at modellene vil komme med de samme svarene. Disse modellene oppdateres ofte og dermed kan metodene modellene bruker til å verdivurdere selskaper og måten de reagerer på input endres.

Vi observerte at modellene gjorde ulik opprinnelig verdivurdering og investeringsbeslutning når vi gjennomførte forsøket ved to anledninger. På grunn av bruksbegrensninger i programvaren vi brukte fikk vi ikke beholdt samme opprinnelige verdivurdering i begge forsøkene. Det måtte gjøres på nytt for begge forsøkene. De opprinnelige vurderingene i forsøk nummer en og to var relativt like, men ikke identiske. Det kan være flere årsaker til at disse ikke var helt like. Det kan skyldes at vi gjennomførte testene med flere dagers mellomrom og at i løpet av den tiden hadde modellene endret måten de fungerte på. Vi kunne heller ikke sjekke om våre tidligere interaksjoner med modellene kunne ha hatt en påvirkning på hvordan modellen vurderte selskapene i fremtiden. Vi fant også noe forskjell i hvordan modellene reagerte på ulik ordlyd, men om denne forskjellen skyldes ulik ordlyd eller at vi testet på forskjellige dager er uvisst. Derfor kan vår reliabilitet være noe svekket. Vi tror selv om at det hadde liten påvirkning, da det ikke var noen nye store offisielle oppdateringer av modellene i mellomtiden av disse to forsøkene.

Da vi testet hvordan modellene reagerte på ulik ordlyd begrenset vi oss til å bruke to ulike ordlyder. Den første ordlyden var av den mer bestemte typen og kunne oppfattes som at vi krevde en ny verdivurdering, mens den andre ordlyden kunne oppfattes som å gi modellene mer frie tøyler. Begrensningene ved vår test av ordlyd var at vi ikke fikk testet dette i et ekte senario. I den virkelige verden ordlegger ikke alle eksperter seg likt, heller ikke ikke-eksperter. Noen ville gitt lange og upresise input, mens andre er kortfattet og konkrete. Derfor kunne det være interessant å teste

hvordan ulik ordlyd kan påvirke modellene ved å gjennomføre en undersøkelse med folk både fra fagmiljøet og amatører.

Verdivurderingene modellene gjorde av de ulike startup selskapene har vi ikke sjekket om stemmer med faktisk verdi på selskapene. Vi ønsket å gjøre en verdivurdering av dette i tillegg, men på grunn av at det er en stor jobb å verdivurdere 50 startups i tillegg til det arbeidet vi har lagt ned i denne oppgaven hadde vi dessverre ikke mulighet til å gjennomføre dette selv. Derfor hadde det var veldig spennende om noen ville gjort en verdivurdering av de selskapene vi har brukt i denne studien og sjekket dette opp mot hva de ulike modellene vurderte selskapene til. Da kunne man sjekket om LLMs er god på å estimere riktig verdi på startups eller ikke.

## 6. Konklusjon

Denne studien har undersøkt hvordan menneskelig input, spesifikt fra eksperter og ikke-eksperter påvirket store språkmodellers verdivurdering av oppstartsselskaper. Studien har benyttet litteratur på effektivitet til LLMs i verdivurdering (Farahani, 2024; Bonaparte, 2024), venture capital kriterier (Gompers et al., 2016; Pintado et al., 2007; Kaplan & Strömberg, 2000), kvalitet på menneskelig input (Schneider, 2020; Amershi et al., 2019), beslutningsbias (Kahneman et al., 1982), og samarbeid mellom mennesker og KI (Mollick, 2024). Med dette teoretiske grunnlaget søkte studien å gi helhetlig bilde av hvordan menneskelig input i form av ekspert- og ikke-ekspertinput påvirket språkmodellene og utlede hypotesene.

For å teste hypotesene valgte vi å gjøre kjikvadrattester, logistisk regresjonsanalyse, t-tester og lineære regresjonsanalyser. Dette førte til resultatene som viste at det ikke var en signifikant forskjell i hvordan språkmodellene, helhetlig responderte til input fra ekspert mot ikke-ekspert når det gjaldt endring av investeringsbeslutning.

Likevel, viste det seg at Gemini sin modell skilte på ekspert- og ikke-ekspertinput. Modellen la mer vekt på ekspertinput og reduserte verdivurderingen av startup selskapene, spesielt når den opprinnelige vurderingen var positiv. Dette tydet på at denne modellen kunne være mer sensitiv til negativ feedback fra eksperter, selv om modellene ikke nødvendigvis endret sin overordnede investeringsbeslutning.

Studien avslørte også signifikante forskjeller mellom Gemini og de andre språkmodellene. Gemini, ut ifra våre funn, ble mer påvirket av ekspertinput enn de andre modellene og var mer sannsynlig til å justere sine investeringsbeslutninger basert på dette. Overordnet hadde Gemini kun en observasjon der den endret en investeringsbeslutning fra negativ til positiv. Disse funnene understreker behovet for å være oppmerksomme til de spesifikke karakteristikene og mulige bias i de ulike språkmodellene når disse benyttes i beslutningsstøtte.

Studien har implikasjoner for bruk av språkmodeller i Venture Capital-bransjen. Den foreslår at språkmodeller kan være verdifulle verktøy som støtte i beslutningsprosesser, men at det er svært

viktig å kritisk evaluere deres output og unngå å stole blindt på deres anbefalinger. Menneskelig dømmekraft og ekspertise forblir essensielt for å gjøre informerte investeringsbeslutninger. Videre forskning bør sette søkelys på bedre forståelse av hvordan språkmodellene integrerer menneskelig input i deres evalueringer og hvordan man kan optimalisere prompts for å oppnå mer nøyaktige resultater.

## 7. KI Erklæring

Erklæring i forbindelse med hvordan kunstig intelligens i studien og hvilke modeller som er benyttet.

VERKTØY	FORMÅL VED BRUK AV VERKTØY
ChatGPT 4o	- Datainnsamling til studien - Sparrepartner og tips til å forbedre tekst - Bistand koding i RStudio
Llama-3.1 405B	- Datainnsamling til studien
Gemini 1.5 Pro	- Datainnsamling til studien
Claude 3.5 Sonnet	- Datainnsamling til studien
Elicit: The AI Research Assistant	- Litteratursøk
Perplexity.ai	- Litteratursøk

*Vi er klar over at vi er ansvarlig for alt innhold i denne masteroppgaven, inkludert de deler der KI-verktøy er benyttet. Vi har ansvar for at oppgaven følger etiske regler for personvern og publisering.*

## 8. Referanser

Affinity. (2024). *AI tools for Venture Capitalists: Enhancing deal flow and decision-making*.

Hentet 21.11.24 fra <https://www.affinity.co/guides/vc-ai-tools>

Ajiga, D.I., Adeleye, R.A., Asuzu, O.F., Owolabi, O.R., Bello, B.G., & Ndubuisi, N.L. (2024).

*Review of AI techniques in financial forecasting: Applications in stock market analysis*.

Finance & Accounting Research Journal.

Akkaya, M. (2020). *Startup Valuation*. Advances in Business Information Systems and Analytics.

Amershi, S., Weld, D.S., Vorvoreanu, M., Fourney, A., Nushi, B., Collisson, P., Suh, J., Iqbal, S.T., Bennett, P.N., Quinn, K.I., Teevan, J., Kikin-Gil, R., & Horvitz, E. (2019).

*Guidelines for Human-AI Interaction*. Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems.

Anthropic. (2024). *Claude 3.5 Sonnet (24. september til 16. oktober versjon) [Stor språkmodell]*.

<https://poe.com/>

Bonaparte, Y. (2024). *Artificial Intelligence in Finance: Valuations and Opportunities*. Finance Research Letters.

Bronzini, M., Nicolini, C., Lepri, B., Passerini, A., & Staiano, J. (2023). *Glitter or gold?*

*Deriving structured insights from sustainability reports via large language models*. EPJ Data Science, 13, 1- 41.

Budiu, R., & Moran, K. (2021, 25. juli). *How Many Participants for Quantitative Usability Studies: A Summary of Sample-Size Recommendations*. NN/g.

<https://www.nngroup.com/articles/summary-quant-sample-sizes/>

Dahl, M. (2004). *Født til succes? Knopskydninger som kilde til innovation*. Ledelse & Erhvervsøkonomi, 68, 285-292.

- 
- Dahlum, S. (2023, 14. januar). *Uavhengig variabel (forskningsmetode)*. I Store norske leksikon på snl.no. [https://snl.no/uavhengig\\_variabel\\_-\\_forskningsmetode](https://snl.no/uavhengig_variabel_-_forskningsmetode)
- Dahlum, S. (2024, 11. mars). *Avhengig variabel (forskningsmetode)*. I Store norske leksikon på snl.no. [https://snl.no/avhengig\\_variabel\\_-\\_forskningsmetode](https://snl.no/avhengig_variabel_-_forskningsmetode)
- DATAtab. (u.å). *Chi-square test*. Hentet 8. desember, 2024, fra <https://datatab.net/tutorial/chi-square-test>
- Dealflow. (u.å). *Vi kobler investorer med selskaper*. Hentet 24. september, 2024–30. september, 2024 fra <https://dealflow.no/>
- De nasjonale forskningsetiske komiteene. (2024). *Etiske dilemmaer ved bruk av kunstig intelligens i forskning*. Hentet 16.11.24 fra <https://www.forskningsetikk.no/>
- Dellermann, D., Lipusch, N., Ebel, P.A., Popp, K.M., & Leimeister, J.M. (2017). *Finding the Unicorn: Predicting Early Stage Startup Success Through a Hybrid Intelligence Method*. *Econometric Modeling: Corporate Finance & Governance eJournal*.
- Elastic. (u.å.a). *What are large language models (LLMs)?* Hentet 12. september, 2024 fra <https://www.elastic.co/what-is/large-language-models>
- Elastic. (u.å.b). *What is natural language processing (NLP)?* Hentet 12. september, 2024 fra <https://www.elastic.co/what-is/natural-language-processing>
- Farahani, M.S. (2024). *Analysis of business valuation models with AI emphasis*. *Sustainable Economies*, 2(3). <https://doi.org/10.62617/se.v2i3.132>
- Felten, Edward W. and Raj, Manav and Seamans, Robert. (April 10, 2023). *Occupational Heterogeneity in Exposure to Generative AI*. <http://dx.doi.org/10.2139/ssrn.4414065>
- Fernández, P. (2002). *Company Valuation Methods*. Accounting Educator: Courses.

- Franke, N., Gruber, M., Harhoff, D., & Henkel, J. (2008). *Venture Capitalists' evaluations of start-up teams: Trade-offs, knock-out criteria, and the impact of VC experience*. *Entrepreneurship Theory and Practice*, 32(3), 459-483. <https://doi.org/10.1111/j.1540-6520.2008.00236.x>
- Galileo. (u.å). *Unlock Evaluation Intelligence for AI Teams*. Hentet 28. november, 2024 fra <https://www.galileo.ai/>
- Galileo. (2024). *LLM Hallucination Index RAG SPECIAL*. <https://galileo.ai/ty/hallucinationindex>
- Gompers, P.A., Gornall, W., Kaplan, S.N., & Strebulaev, I.A. (2016). *How Do Venture Capitalists Make Decisions?* Ewing Marion Kauffman Foundation Research Paper Series.
- Google. (2024). *Gemini 1.5 Pro128k (24. September til 16. Oktober versjon) [Stor språkmodell]*. <https://poe.com/>
- Google. (u.å). *What are AI hallucinations?* Hentet 30. oktober, 2024 fra <https://cloud.google.com/discover/what-are-ai-hallucinations>
- Griffin, L., Kleinberg, B., Mozes, M., Mai, K., Vau, M.D.M., Caldwell, M., & Mavor-Parker, A. (2023). *Large language models respond to influence like humans*. *Proceedings of the First Workshop on Social Influence in Conversations (SICon 2023)*, 15–24. <https://aclanthology.org/2023.sicon-1.3>
- Grønmo, S. (2024, 24. juni). *Validitet*. I Store norske leksikon på [snl.no](https://snl.no/validitet). <https://snl.no/validitet>
- Guinness, H. (2024, 05. august). *The best large language models (LLMs) in 2024*. <https://zapier.com/blog/best-llm/>
- Hu, L. (2024). *Current Study and Method on Artificial Intelligent-based on Venture Capital Decision*. *Proceedings of the 4th International Conference on Economic Management and Big Data Applications, ICEMBDA 2023, October 27–29, 2023, Tianjin, China*.



---

IngestAI. (2024, 20. juni). *How AI improves Venture Capital decision-making.*

<https://ingestai.io/blog/ai-improves-vc-decision-making>

Johannesen, M. (Programleder). (2021–nåtid). #pengepodden – Finn de beste investeringscasene i tidlig fase [Audiopodkast]. Nordnet. <https://www.nordnet.no/blogg/pengepodden-finn-de-beste-investeringscasene-i-tidlig-fase/>

Johannesen, M. (2021, 12. juli). *Hvordan verdsette et tidligfase-selskap?*

<https://www.nordnet.no/blogg/hvordan-verdsette-et-tidligfase-selskap/>

Kahneman, D., Slovic, P., & Tversky, A. (1982). *Judgment Under Uncertainty: Heuristics and Biases*. Cambridge University Press.

Kalinsky, E., & Evtushenko, V. (2023). *Features of valuation of startup projects in entrepreneurial activity*. Herald of Khmelnytskyi National University. Economic sciences.

Kaplan, S.N., & Strömberg, P. (2000). *How Do Venture Capitalists Choose Investments?*

Keeley, R.H., Punjabi, S.S., & Turki, L.A. (1996). *Valuation of Early-Stage Ventures: Option Valuation Models vs. Traditional Approaches*. The Journal of Entrepreneurial Finance. <https://doi.org/10.57229/2373-1761.1186>

Kerner, S.M. (2024, mai) *What are large language models (LLMs)?*

<https://www.techtarget.com/whatis/definition/large-language-model-LLM>

Kostya. (2024, mai) *LLM Bias: Understanding, Mitigating and Testing the Bias in Large Language Models*. <https://academy.test.io/en/articles/9227500-llm-bias-understanding-mitigating-and-testing-the-bias-in-large-language-models>

Liquidity. (u.å.). *How AI can automate loan application approvals and lending*. Hentet

11.12.2024 fra <https://www.liquidity.com/resource-funding/how-ai-can-automate-loan-application-approvals-and-lending>

- Luo, X., Dafflon, J., Bao, D.D., & Love, B.C. (2024). *Large language models surpass human experts in predicting neuroscience results*. arXiv preprint arXiv:2403.03230. <https://arxiv.org/abs/2403.03230>
- Meta. (2024). *Llama 3.1 405B (24. september til 16. oktober versjon) [Stor språkmodell]*. <https://poe.com/>
- Microsoft. (2024, 20 mai). *Understand tokens*. <https://learn.microsoft.com/en-us/dotnet/ai/conceptual/understanding-tokens>
- Miettinen, M.R., & Niskanen, M. (2015). *Lender evaluations of start-up business prospects*. *Managerial Finance*, 41, 102-120. <https://doi.org/10.1108/MF-10-2013-0284>
- Mollick, E. (2024). *Co-intelligence: Living and working with AI*. Portfolio.
- Montani, D., Gervasio, D., & Pulcini, A. (2020). *Startup Company Valuation: The State of Art and Future Trends*. *International Business Research*, 13, 31-45.
- OpenAI. (2024). *ChatGPT 4o (24. september til 16. oktober versjon) [Stor språkmodell]*. <https://chat.openai.com/>
- OpenAI. (2024). *ChatGPT 4o (3. oktober til 16. oktober versjon) [Stor språkmodell]*. <https://poe.com/>
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C.L., Mishkin, P. & Lowe, R. (2022). *Training language models to follow instructions with human feedback*. arXiv preprint arXiv:2203.02155. <https://arxiv.org/abs/2203.02155>
- Pandey, M.K., & Sergeeva, I. (2022). *Artificial Intelligence Impact Evaluation: Transforming Paradigms in Financial Institutions*. *World of Economics and Management*.

- 
- Pew Research Center. (2019, 2. august). *Americans often trust practitioners more than researchers, but are skeptical about scientific integrity*. Pew Research Center. <https://www.pewresearch.org/science/2019/08/02/americans-often-trust-practitioners-more-than-researchers-but-are-skeptical-about-scientific-integrity/>
- Pintado, T.R., De Lema, D.G., & Van Auken, H.E. (2007). *Venture Capital in Spain by Stage of Development*. *Journal of Small Business Management*, 45, 68-88.
- Saunders, M.N., Lewis, P., & Thornhill, A. (2019). *Research Methods for Business Students*. Pearson.
- Schneider, J. (2020). *Humans learn too: Better Human-AI Interaction using Optimized Human Inputs*. arXiv, abs/2009.09266.
- Setlur, V., & Birnbaum, L. (2024). *Can Nuanced Language Lead to More Actionable Insights? Exploring the Role of Generative AI in Analytical Narrative Structure*. arXiv, abs/2405.02763.
- Sharma, N., Liao, Q.V., & Xiao, Z. (2024). *Generative Echo Chamber? Effects of LLM-Powered Search Systems on Diverse Information Seeking*. arXiv, abs/2402.05880.
- Skovlund, E. (2017). Når bør man velge en ikke-parametrisk metode? *Tidsskrift for Den norske legeforening*. Doi: 10.4045/tidsskr.17.0219
- Statistics Solutions. (2024). *Assumptions of Logistic Regression*. Lest 06.11.24 fra <https://www.statisticssolutions.com/free-resources/directory-of-statistical-analyses/assumptions-of-logistic-regression/>
- Sullivan, P.H. (2000). *Valuing intangibles companies – An intellectual capital approach*. *Journal of Intellectual Capital*, 1, 328-340.
- Svartdal, F. (2020, 3. april). *Reliabilitet*. I Store norske leksikon på snl.no. <https://snl.no/reliabilitet>

Talboy, A.N., & Fuller, E. (2023). Challenging the appearance of machine intelligence: Cognitive bias in LLMs. *ArXiv*, [abs/2304.01358](https://arxiv.org/abs/2304.01358).

Tan, W.C. (2023). *Unstructured and structured data: Can we have the best of both worlds with large language models?* *IEEE Data Eng. Bull.*, 46, 5-11.

Tuli, Mantej Singh., Shreyas Gav, Dheemanth A Na, Sree Chand Ra, Anupama Y Ka (2023). *Smart Start-Up Analyzer: Prediction Model, Analysis tool for Venture Capitals using Machine Learning*. *International Journal of Scientific Research in Engineering and Management*.

University of Michigan-Dearborn. (2023). *AI's mysterious 'black box' problem, explained*. <https://umdearborn.edu/news/ais-mysterious-black-box-problem-explained>

View Group. (2023, 17. mars). *Hvordan verdsette et selskap?* <https://viewgroup.no/verdsette-et-selskap/>

Visible. (2024, 17. april). *AI tools for Venture Capital: Revolutionizing insights and decision-making*. <https://visible.vc/blog/ai-tools-for-venture-capital>

Wilimowska, Z., & Krzysztozek, T. (2013). *The Use of Artificial Neural Networks in Company Valuation Process*. In N. Nguyen, B. Trawiński, R. Katarzyniak, & G.S. Jo (Eds.), *Advanced Methods for Computational Collective Intelligence* (Vol. 457). Springer. [https://doi.org/10.1007/978-3-642-34300-1\\_27](https://doi.org/10.1007/978-3-642-34300-1_27)

Виноградов, А. (2023). *Methodology for assessing the market value of Venture Capital companies*. *Вестник Академии права и управления*. [https://doi.org/10.47629/2074-9201\\_2023\\_2\\_128\\_135](https://doi.org/10.47629/2074-9201_2023_2_128_135)