

KI og opphavsrett

Tilstandsrapport våren 2024

Innledning

Generativ kunstig intelligens og opptrening av språkmodeller er et område i stadig endring. Politikk- og strategiutforming som er fundert på en solid og oppdatert kunnskapsplattform er avgjørende for å få gjennomslag i et landskap der aktørene er mektige og pågående.

Kopinors KI©K-prosjekt startet opp i desember 2023 og har som mål å gjøre administrasjonen, styret og medlemsorganisasjonene bedre i stand til å utforme politikk og strategi og bedre kunne vurdere risiko og mulig gevinst ved fremtidig forretningsutvikling, spesielt knyttet til å tillate og lisensiere bruk av verk i språkmodeller og maskinlæring.

Arbeidet vil være preget av at KI-feltet er under utvikling i et høyt tempo. Prosjektet vil særlig vektlegge juridiske og teknologiske faktorer i utredningen. I tillegg må rettighetshavernes interesser inkluderes i vurderingene.

Denne første rapporten tar for seg utviklingen og noen av utfordringene innenfor lovgivning og jus. Feltet er i flyt og utvikling i et ganske høyt tempo, og dette er ikke tiden for definitive konklusjoner. Målet er i stedet å legge grunnlaget for en diskusjon som flere kan delta i, med håp om et felles, norsk ståsted eller i det minste opplyst uenighet.

Vi har i denne rapporten valgt å konsentrere oss om innhold i form av tekst. Tilsvarende problemstillinger gjelder imidlertid også musikk, visuell kunst og annet grafisk utformet innhold, som også inngår i tjenester som bygger på generativ kunstig intelligens.

Rapporten inneholder en del ord og begreper knyttet til KI. Mange av disse begrepene er forsøkt forklart så enkelt som mulig i slutten av rapporten. Der ligger det også en tidslinje som anskueliggjør hvor raskt både teknologiutvikling og juridiske dilemmaer har beveget seg siden lanseringen av ChatGPT rett før jul i 2022.

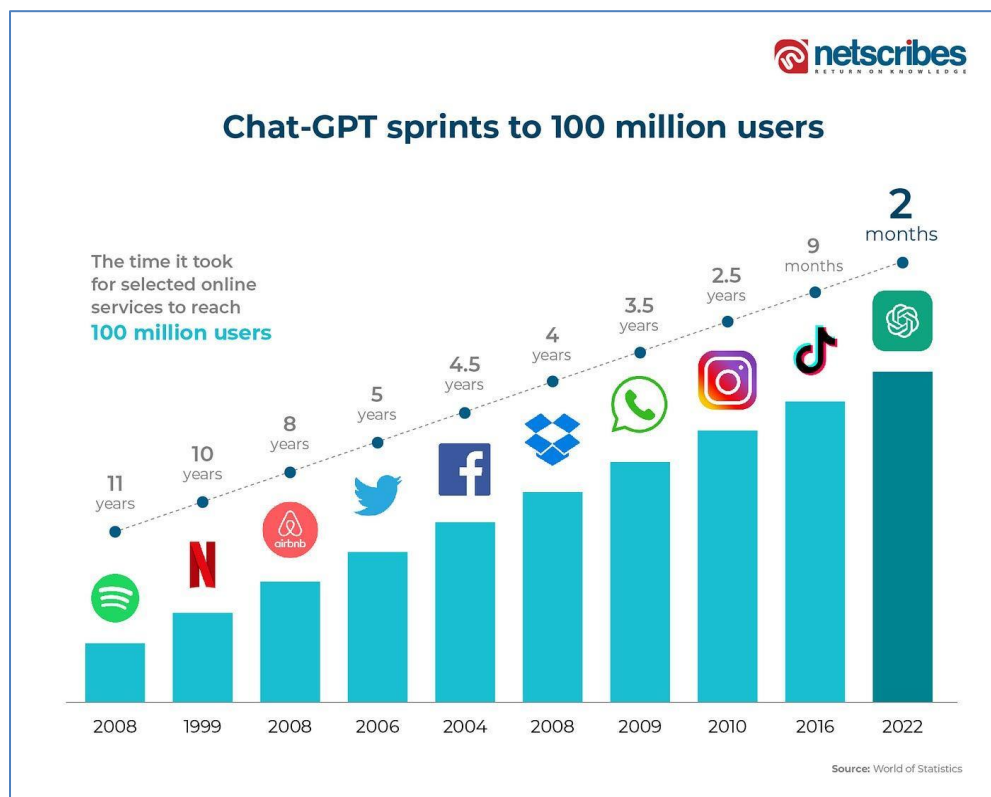
Rapporten er skrevet av Geir Terje Ruud, i samarbeid med Frank Gander.

Innholdsfortegnelse

Innledning.....	1
1. KI-revolusjonen.....	3
Søksmål.....	4
Juridiske tilbakeslag	6
Nye partnerskap.....	7
2. Juridiske problemstillinger.....	8
EUs digitalmarkedsdirektiv.....	8
EUs AI Act.....	9
Er maskinlæring tekst- og datautvinning?.....	10
Opt-out eller opt-in	11
«Fair use» av bøker og artikler	13
Nordiske tiltak	13
3. Nasjonalbibliotekets rolle	14
Prosjekt Mímir	14
Rettigheter.....	15
Kompensasjonsordning	16
Tillegg.....	17
Tidslinje	17
Ordlister.....	19

1. KI-revolusjonen

Open AIs ChatGPT kom som en overraskelse på de fleste senhøsten 2022. Aldri før har ny teknologi blitt tatt i bruk raskere av store grupper. Søylen under viser hvor få dager det gikk fra lansering til 100 millioner hadde tatt ChatGPT i bruk. Søylene som viser hvordan Instagram, TikTok, Netflix, Spotify og andre teknologiske nyvinninger og suksesshistorier vokste frem, dokumenterer hvor voldsom ChatGPTs vekst var.



Men dette var ikke første gang vi tok KI i bruk. Mange av oss brukte for eksempel allerede KI til å sortere bilder på telefonen, med ansiktsgjenkjenning.

Det var heller ikke første gang språkmodeller ble brukt til avanserte oppgaver. Mindre språkmodeller hadde tidligere blitt brukt til for eksempel oversettelse. NTB (Norsk Telegrambyrå) bygde en «nynorsk-robot» på grunnlag av rundt 30.000 tekster som var manuelt oversatt fra bokmål til nynorsk av byråets medarbeidere. Ved hjelp av maskinlæring og manuell innsats gikk oversetteren fra å være rundt 70 prosent av perfekt til et nivå der suksessraten er nesten 98 prosent.

Det revolusjonerende med ChatGPT var den gode dialogfunksjonen. Når vi brukere stiller et spørsmål, svarer den med gode setninger som svært ofte gir mening. Men setningene bygger først og fremst på sannsynlighet, ikke på kunnskap. ChatGPT bruker sannsynlighetsberegning for å finne ut hvilket ord som burde etterfølge det forrige. Jo bedre du forklarer konteksten, jo større er sjansen for et godt svar.

Språkmodellene er fremstilt ved å bli trent på store mengder tekst. Hvilke muligheter som finnes for at innhold skal kunne «trekkes ut» fra modellen, er uklart. Disse og andre tekniske utfordringer vil vi komme tilbake til i en senere rapport om teknologi.

Søksmål

ChatGPT fra selskapet OpenAI er bygget ved at milliarder av setninger er matet inn i en gigantisk database. Nesten 100 språk er representert i den digitale assistent-tjenesten, deriblant norsk. Det er ikke søkt om tillatelse til å bruke innholdet som ligger til grunn for språkmodellene, og det er uklart hvor stor del av modellene som er skapt på grunnlag av opphavsrettsbeskyttet materiale, f.eks. bøker.

Maskinlæringen skjer på grunnlag av en rekke store datasett, blant annet innhold «skrapet» fra nettet. Det er kjent at to store bok-korpus (kalt *books1* og *books2*) inngår, og disse skal visstnok basere seg på bøker som er falt i det fri eller som av andre grunner ikke er opphavsrettsbeskyttet. Imidlertid finnes det et datasett av piratkopierte bøker, *books3*, som består av ca. 200 000 bøker. Meta har erkjent at dette datasettet er brukt i treningen av deres Llama-modeller.

Oversikten under fra blant andre PhD-studenten Kent K. Chang på Berkeley i California¹, viser litt om hvor stor del av noen bøker som er gjenkjent i henholdsvis dialogtjenesten ChatGPT, språkmodellen GPT-4 og Googles BERT. 0.98 betyr at 98 prosent av en bok er gjenkjent.

Treningsgrunnlaget for BERT er i utgangspunktet kjent og skal bestå av Wikipedia og et mindre bok-korpus hentet fra selvpubliserende forfattere (men som overraskende nok også omfatter *Fifty Shades of Grey*).

GPT-4	ChatGPT	BERT	Date	Author	Title
0.98	0.82	0.00	1865	Lewis Carroll	<i>Alice's Adventures in Wonderland</i>
0.76	0.43	0.00	1997	J.K. Rowling	<i>Harry Potter and the Sorcerer's Stone</i>
0.74	0.29	0.00	1850	Nathaniel Hawthorne	<i>The Scarlet Letter</i>
0.72	0.11	0.00	1892	Arthur Conan Doyle	<i>The Adventures of Sherlock Holmes</i>
0.70	0.10	0.00	1815	Jane Austen	<i>Emma</i>
0.65	0.19	0.00	1823	Mary W. Shelley	<i>Frankenstein</i>
0.62	0.13	0.00	1813	Jane Austen	<i>Pride and Prejudice</i>
0.61	0.35	0.00	1884	Mark Twain	<i>Adventures of Huckleberry Finn</i>
0.61	0.30	0.00	1853	Herman Melville	<i>Bartleby, the Scrivener</i>
0.61	0.08	0.00	1897	Bram Stoker	<i>Dracula</i>
0.61	0.18	0.00	1838	Charles Dickens	<i>Oliver Twist</i>
0.59	0.13	0.00	1902	Arthur Conan Doyle	<i>The Hound of the Baskervilles</i>
0.59	0.22	0.00	1851	Herman Melville	<i>Moby Dick; Or, The Whale</i>
0.58	0.35	0.00	1876	Mark Twain	<i>The Adventures of Tom Sawyer</i>
0.57	0.30	0.00	1949	George Orwell	<i>1984</i>
0.54	0.10	0.00	1908	L. M. Montgomery	<i>Anne of Green Gables</i>
0.51	0.20	0.01	1954	J.R.R. Tolkien	<i>The Fellowship of the Ring</i>
0.49	0.16	0.13	2012	E.L. James	<i>Fifty Shades of Grey</i>
0.49	0.24	0.01	1911	Frances H. Burnett	<i>The Secret Garden</i>
0.49	0.12	0.00	1883	Robert L. Stevenson	<i>Treasure Island</i>
0.49	0.16	0.00	1847	Charlotte Brontë	<i>Jane Eyre: An Autobiography</i>
0.49	0.22	0.00	1903	Jack London	<i>The Call of the Wild</i>

Table 1: Top 20 books by GPT-4 name cloze accuracy.

¹ [GPT-4 memorizes contents of copyrighted books, and it could be a cultural issue \(the-decoder.com\)](https://the-decoder.com)

The Atlantic lagde høsten 2023 en løsning hvor man kunne søke etter sin egen utgivelse i datasettet *books3*. Løsningen viser at også en rekke tekster av norske forfattere har blitt brukt i trening av språkmodeller.

- [Her kan du teste The Atlantics søk i datasettet *books3*](#)



Foto: NTB / Claudio Bresciani

Norske bøker brukt til å trene AI-roboter: – Det er fullstendig uakseptabelt

En rekke norske forfatters bøker er brukt til å trene opp flere store AI-roboter – uten forfatterens tillatelse.

VG, 27. september 2023

Det er sterke grunner til å tro at det er svært mye innhold fra redaktørstyrte medier i modellene som nå er tilgjengelig. News Media Alliance skriver blant annet i sin gjennomgang:

«News and media content is overrepresented in samples of popular curated sets such as C4, OpenWebText, or OpenWebText2 used for LLM training, as compared to the broader category of material captured in the Common Crawl. [...]

News and media reports ranks third among all categories of sources in Google's C4 data set, and half of the top ten represented sites overall are news outlet.»²

Det er særlig verdt å trekke frem WebText, et av OpenAIs datasett som de selv har omtalt som et «high quality» datasett og som i stor grad består av nyhetsinnhold.

New York Times har gått til søksmål mot OpenAI og deres hovedeier Microsoft blant annet fordi de mener ChatGPT nå er «en konkurrent om å formidle troverdig informasjon». Det har de blitt fordi de – uten lov – har brukt New York Times-artikler til å bygge ChatGPT, ifølge søksmålet fra 27. desember 2023.

New York Times skriver blant annet at OpenAI har brukt alt deres materiale og dermed tilegnet seg ikke bare setninger og ord, men også verdifull kunnskap som bygger verdi for Open AI, ikke for New York Times. De dokumenterer at OpenAI gjenskaper New York Times-artikler med overraskende presisjon. Dette punktet er interessant nettopp fordi det viser at generative KI-tjenester kan konkurrere på samme markeder som originaltekstene.

Både OpenAI og Microsoft har i mars 2024 lagt fram saksframlegg der de ber retten avvise søksmålet. Når artikler blir kopiert ordrett, bygger det på feil i modellene, skriver OpenAI.

Også amerikanske Authors Guild har saksøkt OpenAI på vegne av sine 14 000 medlemmer. Jonathan Franzen har blitt brukt som «poster boy» for søksmålet og sier:

«Generative AI is a vast new field for Silicon Valley's longstanding exploitation of content providers. Authors should have the right to decide when their works are used to 'train' AI. If they choose to opt in, they should be appropriately compensated.»³

På billedsiden ble KI-selskapene Midjourney, Deviantart og Stability AI allerede i januar 2023 saksøkt av en rekke amerikanske billedkunstnere i et gruppesøksmål. Senere gikk billedbyrået Getty Images til sak mot Stability AI for bruken av deres bilder i treningsmaterialet.

I en rapport fra plagiat-tjenesten Copyleaks, publisert i slutten av februar 2024, hevdes det at så mye som 60 prosent av svarene fra Open AIs 3.5-språkmodell inneholder en eller annen form for plagiert innhold.⁴

Juridiske tilbakeslag

Et eget søksmål, der blant annet Sarah Silverman er sentral, fikk midt i februar avvist deler av kravene sine fra en dommer.⁵ Dommeren mener for eksempel at det ikke er sannsynliggjort at man kan hente ut hele verk fra OpenAIs løsning og ser det som «for spekulativt å hevde at fremtidig verdi av materialet er påført stor skade». Men dommeren anerkjenner at verkene er

² [White Paper: How the pervasive copying of expressive works to train and fuel generative artificial intelligence systems is copyright infringement and not a fair use \(News Media Alliance\)](#)

³ [The Authors Guild, John Grisham, Jodi Picoult, David Baldacci, George R.R. Martin, and 13 Other Authors File Class-Action Suit Against OpenAI \(The Authors Guild\)](#)

⁴ [New report: 60% of OpenAI model's responses contain plagiarism \(axios.com\)](#)

⁵ [Sarah Silverman Hits Stumbling Block in AI Lawsuit Against Meta \(hollywoodreporter.com\)](#)

brukt til å skape språkmodellen, uten å spørre Silverman om lov. Saksøkerne har frist til 15. mars for å justere sitt søksmål.

Også sakene mot Midjourney og Deviantart har møtt tilsvarende delvise avvisninger i rettsapparatet. Det vil trolig ta minst tre år i amerikanske rettsinstanser før det foreligger en rettskraftig dom i disse sakene.

Nye partnerskap

I motsetning til New York Times, har det tyske Axel Springer-konsernet valgt å inngå en avtale med OpenAI. Dette partnerskapet er det første av sitt slag mellom et forlag/mediehus og et KI-forskningselskap.

Avtalen innebærer at Axel Springer vil lisensiere aktuelt nyhetsinnhold og arkivmateriale til OpenAI for bruk i trening av OpenAIs språkmodeller, som ChatGPT. Axel Springer, som utgir publikasjoner som Business Insider og Politico, vil tillate ChatGPT å oppsummere aktuelle artikler i sine svar.

ChatGPT-brukere over hele verden vil motta sammendrag av utvalgt globalt nyhetsinnhold fra Axel Springer, inkludert materiale som ellers ville vært bak en betalingsmur, og med lenker til de fullstendige artiklene. Axel Springer vil bli kompensert for å gjøre sitt innhold tilgjengelig for OpenAI, og avtalen er gyldig i flere år uten å binde noen av partene til eksklusivitet, noe som betyr at de er frie til å inngå nye avtaler med andre parter. Det er ikke kjent hva avtalen betyr økonomisk sett for Axel Springer.

Nyhetsbyrået AP har inngått en liknende avtale med OpenAI. AP bekrefter at samarbeidet med OpenAI inkluderer en økonomisk godtgjørelse ved siden av tilgang til teknologi, men de vil ikke si noe om størrelsen på godtgjørelsen. De gir ikke OpenAI adgang til «current news».

Medieselskapene NewsCorp og Thomson Reuters bekrefter også at de er i tenkeboksen, eller er i forhandlinger med noen av de større KI-selskapene om bruk av deres materiale i trening og/eller kommersielle tjenester.

2. Juridiske problemstillinger

Som vi har sett, har utviklingen innenfor generativ KI allerede ført til starten på flere rettsprosesser i USA. Vi kommer tilbake til noen av problemstillingene nedenfor. Men mens angloamerikansk rett særlig bygger på avgjørelser i domstolene, vil utviklingen i europeisk opphavsrett i større grad ta utgangspunkt i lovgivning. Vi skal derfor starte med de viktigste lovprosessene knyttet til KI i Europa.

EUs digitalmarkedsdirektiv

Digitalmarkedsdirektivet (Digital Single Market Directive, DSM) ble endelig vedtatt i 2019 med frist til juni 2021 for medlemslandene til å gjennomføre direktivet. Direktivet ble av EØS-komiteen i desember 2023 besluttet tatt inn i EØS-avtalen og ligger til grunn for forslaget til endringer i åndsverkloven, som Kulturdepartementet sendte på høring i november 2023.⁶

Direktivet omhandler en rekke spørsmål om digital utnyttelse av innhold og bygger videre på EUs opphavsrettsdirektiv fra 2001. Under behandlingen av digitalmarkedsdirektivet var det særlig oppmerksomhet rundt artikkel 15 (om presseutgiveres rettigheter) og artikkel 17 (om nettplattformenes ansvar for brukeropplagte innhold). For spørsmålet om kunstig intelligens er det imidlertid artikkel 3 og 4 om tekst- og datautvinning (*Text and data mining*, TDM) som er interessant:

- Artikkel 3 åpner for at forskningsinstitusjoner og kulturarvinstitusjoner som har lovlig tilgang til et verk, kan bruke verket til tekst- og datautvinning til egne vitenskapelige forskningsformål.
- Artikkel 4 handler om bruk av rettighetsbeskyttet materiale til tekst- og datautvinning på andre områder. Dette forutsetter at rettighetshaveren ikke har nedlagt forbud mot slik bruk.

Mye av kjernen i diskusjonen vil dreie seg om grensene mellom artikkel 3 og artikkel 4. Hvordan går man opp grensen mellom de to artiklene? Hvor slutter forskningen? På hvilken måte kan resultater fra vitenskapelig forskning brukes i videre utvikling? Flere problemstillinger vil trolig melde seg i denne grenseoppgangen.

Kulturminister Lubna Jaffery (Ap) snakket om DSM-direktivet i Europautvalget på Stortinget 1. februar 2024⁷, og sa dette om stridspunktene tre og fire:

«Reguleringen av opphavsrett skal balansere mellom rettighetshavernes interesser på den ene siden og brukerne og samfunnets interesser på den andre. Et eksempel på hvor denne spenningen kommer til syne og hvor det har vært litt debatt også i Norge, er artiklene 3 og 4 om såkalt tekst- og datautvinning. Disse bestemmelsene vil bl.a. kunne få betydning for opptrening av kunstig intelligens. For forskningsformål blir det en vid adgang til å bruke beskyttet materiale, mens det for kommersielle formål kreves at rettighetshaverne ikke har nedlagt forbud. Man kan altså reservere seg. Blant rettighetshavere er det en bekymring for hvordan en slik reservasjonsrett skal fungere i praksis. Jeg

⁶ [Høring - endringer i åndsverkloven mv. \(Regjeringen.no\)](#)

⁷ [Møte i Europautvalget torsdag den 1. februar 2024 kl. 8:30 \(Stortinget\)](#)

regner med at vi kommer til å få belyst disse problemstillingene nærmere i høringen.»

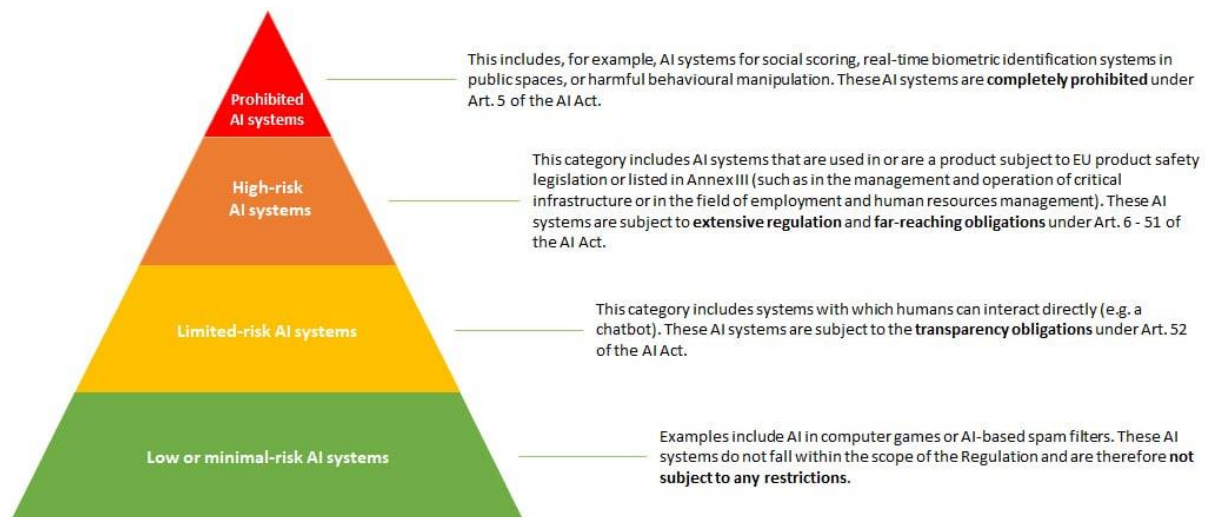
Etter den foreløpige tidsplanen vil endringene i åndsverkloven bli lagt fram for Stortinget våren 2025 og kan tidligst tre i kraft 1. juli 2025.

EUs AI Act

Da medlemslandene i EU fredag 2. februar 2024 sa ja til det eksisterende forslaget til AI Act (norsk: KI-forordningen⁸), var det et gjennombrudd med tanke på fremtidig kontroll og regulering av hva man kan gjøre med KI.

Med AI Act går EU foran i et forsøk på å regulere et område, på samme måte som EU var først med å sikre personvernet gjennom GDPR-forordningen. Det er tydelige krav om blant annet transparens omkring hvordan og om KI er brukt. Dette omfatter også åpenhet om hva slags innhold som er brukt i treningen av generativ KI.

Forpliktelsene i forordningen avhenger av risikoklassifiseringen til det aktuelle KI-systemet. I en oversiktlig og ganske lett forståelig gjennomgang av KI-forordningen som konsultentselskapet Deloitte gjorde før jul 2023, går man spesielt inn på de sikkerhetsmessige sidene og de personvernmessige sidene ved KI.⁹



Oversikt over AI Acts klassifisering av risiko (Deloitte)

Professor Natali Helberger ved Universitetet i Amsterdam er en jevnlig gjest i Oslo, blant annet på OsloMet. I en fersk gjennomgang som hun og en rekke kolleger kaller «*The Amsterdam paper*» trekker de frem hvor mye som kreves med tanke på åpenhet for å følge den nye forordningens regler. Målet er at man med transparens skal kunne sikre rettighetshavere. AI Act krever blant annet at man deklarerer bruk av rettighetsbeskyttet materiale og ikke rettighetsbeskyttet materiale på samme måte.

⁸ [KI-forordningen om europeisk regelverk for kunstig intelligens \(Europalov\)](#)

⁹ [EU AI Act: Hva må du vite og hvordan kan du forberede deg? \(Deloitte\)](#)

Helberger trekker frem at transparens også gjelder for at det skal være forståelig, så det kreves en narrativ forklaring av det som er gjort, ikke kun en teknisk beskrivelse:

«The recital in question makes it clear that the ‘summary should be comprehensive in its scope instead of technically detailed, for example, by listing the main data collections or sets that went into training the model, such as large private or public databases or data archives, and by providing a narrative explanation about other data sources used’. It also emphasises that the template to be provided by the AI Office should ‘allow the provider to provide the required summary in narrative form’.»¹⁰

Åpenhet om risikovurdering blir en viktig del av AI Act, akkurat som med GDPR. Men med risiko tenkes mer på liv og helse, personvern og stor samfunnsmessig påvirkning, heller enn risikoen for å misbruke opphavsrettsbeskyttet materiale. I den endelige versjonen av AI Act er det imidlertid også satt inn direkte henvisninger til opphavsretten i opplistingen av det tilbydere av generative KI-modeller (*general-purpose AI models*) er forpliktet til:

(c) put in place a policy to respect Union copyright law in particular to identify and respect, including through state of the art technologies, the reservations of rights expressed pursuant to Article 4(3) of Directive (EU) 2019/790

(d) draw up and make publicly available a sufficiently detailed summary about the content used for training of the general-purpose AI model, according to a template provided by the AI Office

Arbeidet med KI-forordningen i Norge er lagt til Digitaliserings- og forvaltningsdepartementet.

Er maskinlæring tekst- og datautvinning?

Gjennombruddet for generativ kunstig intelligens inntrådte etter at digitalmarkedsdirektivet ble vedtatt i 2019. Med maskinlæring og generativ KI har altså bruksmåtene for tekst- og datautvinning blitt betydelig utvidet. Det har derfor blitt stilt spørsmål ved om unntakene for tekst- og datautvinning i direktivet virkelig kan omfatte den utstrakte bruken som maskinlæring utgjør.

Foreløpig besvares dette spørsmålet med «ja». Kulturdepartementet er for eksempel i sitt forslag til endret åndsverklov klar på at opptrening av kunstig intelligens omfattes av tekst- og datautvinning. EU knytter også generativ KI direkte til reglene om tekst- og datautvinning i KI-forordningen (se ovenfor)

Likevel vil mange hevde at maskinlæring går lenger enn unntakene for tekst- og datautvinning egentlig var ment for. I Polen, som er det eneste medlemslandet som ikke har gjennomført digitalmarkedsdirektivet, har regjeringen i 2024 lagt fram forslag om at bruk til generativ KI unntas fra bestemmelsene om tekst- og datautvinning.

I Frankrike har medlemmer fra Macrons Renaissance-parti introdusert en endring i loven om immaterielle rettigheter (CPI) for å sikre forholdene mellom kunstig intelligens og

¹⁰ [The Amsterdam Paper: Recommendations for the technical finalisation of the regulation of GPAI in the AI Act \(AI, Media & Democracy Lab\)](#)

opphavsrett.¹¹ Dette lovforslaget introduserer flere nye forpliktelser, men tre av dem er spesielt interessante her:

- Enhver integrering av verk beskyttet av opphavsrett i utgangspunktet er underlagt de generelle bestemmelsene i CPI. Det må derfor gis tillatelse fra opphavspersonen eller rettighetshaveren før verket kan brukes til trening eller tjenester av KI-systemer.
- De legitime rettighetshaverne til et KI-generert verk – skapt uten menneskelig inngripen – tilhører de hvis verk ble brukt til å skape det kunstige verket. I Frankrike er dette sett på som et ganske kontroversielt forslag, og et brudd med eksisterende lovgivning.
- Det legges opp til en ny skatt som skal betales av KI-leverandører til opphavsrettsorganisasjonene for kollektiv forvaltning. Skatten betales for de tilfellene hvor KI-generert innhold er bygget på verk med ukjent opprinnelse.

Lovforslaget er ennå ikke behandlet av den franske nasjonalforsamlingen og må trolig sees opp mot de forpliktelsene og avgrensingene som ligger i EUs direktiver og forordninger.

I Storbritannia har *Intellectual Property Office* opprettet en arbeidsgruppe bestående av rettighetshavere og KI-utviklere for å diskutere samspillet mellom opphavsrett og KI. Arbeidsgruppen har ikke lyktes i å komme til enighet om en effektiv frivillig ordning, Departementet for forskning, innovasjon og teknologi skriver i en rapport at de vil fortsatt engasjere seg i dialogen «for å sikre en levedyktig og effektiv tilnærming som lar KI- og kreative sektorer vokse i partnerskap. Det legges vekt på at KI-utvikling skal støtte, ikke undergrave, menneskelig kreativitet og innovasjon».

Opt-out eller opt-in

Reservasjon mot crawlere er brukt på internett i lang tid, for eksempel ved å gi beskjed om at en nettside ikke skal kunne indekseres av Google. Det kan f.eks. gjøres i kildekoden i filen *robots.txt* ved å følge reglene i standarden *Robots Exclusion Protocol*.

Mange nyhetsorganisasjoner har på samme måte lagt inn sperre for at KI-crawlere skal kunne se deres innhold. Det betyr at en tjeneste som Signal, der en journalist skriver fem-seks setninger og en KI-tjeneste bygger ut artikkelen med informasjon fra andre kilder, ikke har innhold fra for eksempel New York Times. Men de kan ha innhold fra selskap som eies av for eksempel Axel Springer.

Det at innhold fra New York Times ikke kan brukes (fordi de ikke har en avtale om godtgjørelse) betyr, som den danske mediekommentator Jan Birkemose skriver, at de beste og mest troverdige kildene utelates.

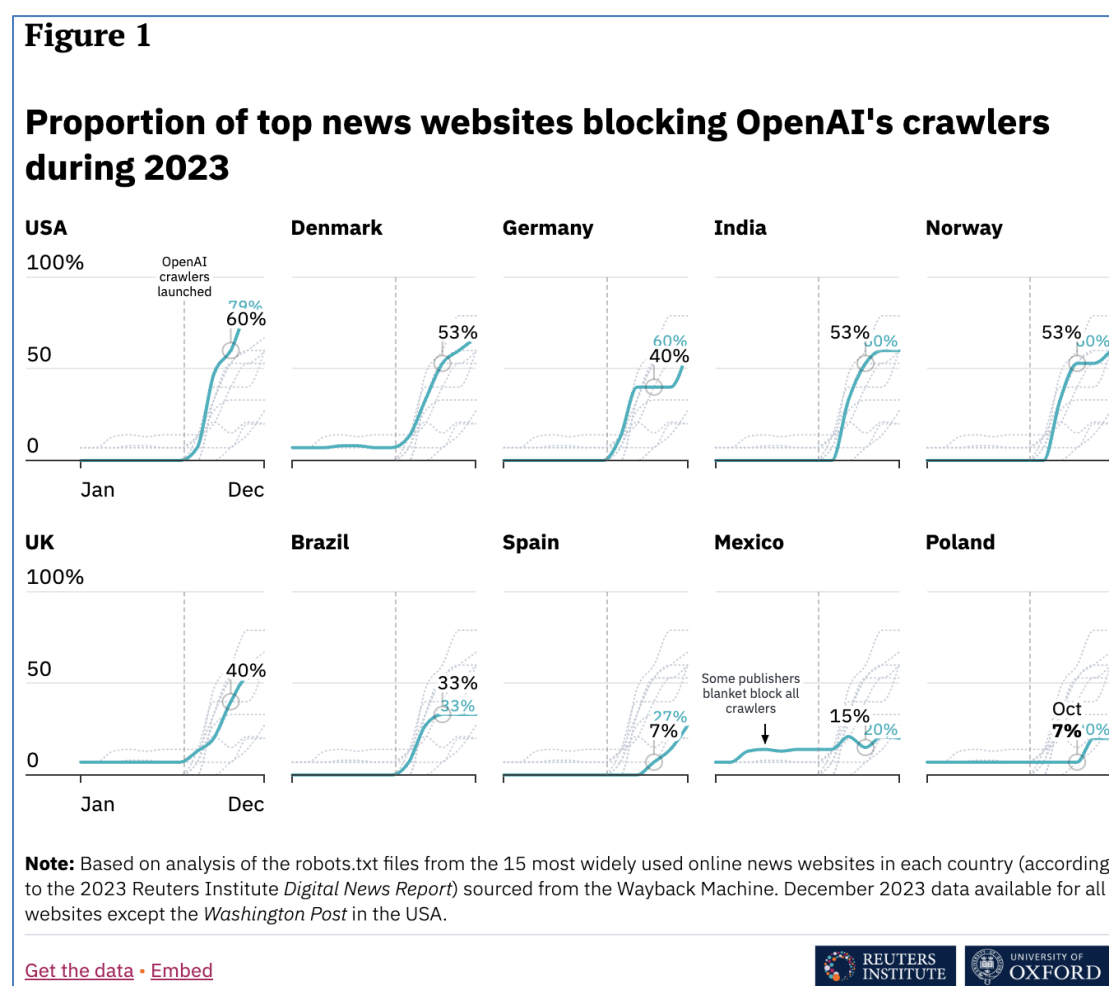
Denne reservasjonsmuligheten omtales ofte som «opt-out». Slik vi viste ovenfor, forutsetter digitalmarkedsdirektivet at verk ikke kan benyttes til tekst- og datautvinning utenom forskningsformål, hvis rettighetshaverne har reservert seg, for eksempel på egen nettside eller ved en maskinlesbar tekst i kildekoden.

¹¹ [French National Assembly Proposes New Law to Secure Copyright in AI \(Fieldfisher\)](#)

Det kan diskuteres om en opt-out-løsning er god nok, eller om man heller burde arbeide for en opt-in-løsning, der man må si ja til å la sitt innhold bli brukt. Dette er prinsippet som gjelder ellers i opphavsretten, men det er altså ikke den løsningen digitalmarkedsdirektivet foreskriver.

Et faktum som kan brukes mot norske rettighetshavere - spesielt fra de internasjonale KI-gigantene - er at man de siste 20 årene stilltiende har latt Google (og andre) indekser alt innhold for å sikre god trafikk tilbake fra søkemotorene. Dette i motsetning til eksempelvis enkelte danske og belgiske rettighetshavere som har nektet Google å gjøre deler av denne indekseringen – spesielt på nyhets- og aktualitetsinnhold.

I en helt fersk rapport fra Reuters Institute for the Study of News på Oxford University¹² dokumenteres det at i Norge har rundt 60 prosent av nyhetsmediene reservert seg mot crawling fra Open AI, mens bare 13 prosent av norske nyhetsmedier har reservert seg mot Googles KI-crawler. Tallene er for 2023.



Det er interessant i et nordisk perspektiv at mens norske og danske nyhetsmedier håndterer reservasjon mot OpenAI noenlunde likt (Danmark: 67 prosent og Norge 60 prosent), så er det 40 prosent av danske nyhetsmedier som reserverer seg mot Googles KI-løsning, mot bare 13 prosent i Norge.

¹² [How many news websites block AI crawlers? \(Reuters Institute for the Study of Journalism\)](#)

«Fair use» av bøker og artikler

KI-selskapene selv mener at trening på opphavsrettslig materiale er svært viktig for å skape gode nok tjenester.

Open AI hevder å overholde gjeldende lover, inkludert opphavsrettslover. Selskapet argumenterer for at trening på opphavsrettsbeskyttet materiale faller inn under «rimelig bruk» (*fair use*), at innholdsskaperes rettigheter respekteres, og at Open AI arbeider aktivt for å finne balanserte løsninger som gagnar både KI-utvikling og opphavsrettsinnehavere.

I et hørings svar til det britiske overhusets digital-komité i desember 2023, understreket Open AI nødvendigheten av omfattende treningsdata, inkludert opphavsrettsbeskyttet materiale, for å utvikle effektive KI-modeller som møter dagens samfunnsbehov:

Because copyright today covers virtually every sort of human expression – including blog posts, photographs, forum posts, scraps of software code, and government documents – it would be impossible to train today’s leading AI models without using copyrighted materials. Limiting training data to public domain books and drawings created more than a century ago might yield an interesting experiment, but would not provide AI systems that meet the needs of today’s citizens.¹³

Modellene kjennetegnes av mangel på transparens, så man kan ikke være helt sikker på hva som er brukt. News Media Alliance har laget et «white paper» som argumenterer for hvorfor slik bruk ikke kan ses på som *fair use*.¹⁴

Nordiske tiltak

I Sverige og Danmark er digitalmarkedsdirektivet allerede gjennomført. Som i Norge arbeides det med å finne en posisjon som svarer både på behovet for å sikre rettigheter og for å kunne delta i teknologiutviklingen.

I Danmark ble det i desember 2023 lagt fram forslag til endringer i opphavsrettsloven, blant annet slik at det kan kreves mekling i tvister mellom tek-selskaper og rettighetshavere. Det ble også satt ned en ekspertgruppe om opphavsrett og KI, presentert på denne måten av kulturminister Jakob Engel-Schmidt:

«Kunstig intelligens kommer grundlæggende til at forandre alle kunstneres arbejdsvilkår og rettigheder. Som kulturminister vil jeg sikre, at de skabende kunstnere ikke ender som tabere, mens techfirmaer som en flok ansvarsløse gribbe spiser af både værker, rettigheder og kreativitet.»¹⁵

Ekspertgruppen skal ifølge pressemeldingen fra kulturministeren levere sin rapport før sommeren.

¹³ [OpenAI—written evidence \(House of Lords Communications and Digital Select Committee Inquiry\)](#)

¹⁴ [White Paper: How the pervasive copying of expressive works to train and fuel generative artificial intelligence systems is copyright infringement and not a fair use \(News Media Alliance\)](#)

¹⁵ [Nu bliver ophavsretsloven moderniseret \(Kulturministeriet\)](#)

3. Nasjonalbibliotekets rolle

Mens forskning og utvikling av generativ KI internasjonalt gjerne skjer i regi av private selskaper, er det her hjemme offentlige institusjoner, særlig universiteter og Nasjonalbiblioteket (NB) som står i første rekke. Spesielt for Norge er Nasjonalbibliotekets sentrale rolle. Dette har sammenheng med bibliotekets offensive digitaliseringsstrategi, der målet er å digitalisere hele samlingen. Eksempelvis er alle bøker digitalisert, og utgivelsene til og med 2005 er lagt ut på nettet i medhold av Bokhylla-avtalen.

I forskrift til åndsverkloven § 4 gis det rom for eksemplarframstilling av verk for forskningsformål. Nasjonalbiblioteket gis her også en spesifikk hjemmel til å fremstille eksemplarer av åndsverk i sine samlinger som grunnlag for språklige korpuser, også kun for forskningsformål. Denne hjemmelen omfatter også pliktavlevert materiale og krever ingen godkjenning.

Gjeldende lov og forskrift er formulert og vedtatt før generativ KI. Nasjonalbiblioteket selv er imidlertid klar på at det finnes begrensninger, slik leder for NBs AI-lab, Svein Arne Brygfjeld, skriver i Stat & Styring 3/2023:

«Tilgangen til NBs samling som treningsdata har dog sine begrensninger av opphavsrettslige grunner. For generell ekstern bruk kan NB bare gi tilgang til innhold som enten ikke er begrenset av opphavsrett eller der det finnes avtaler om bruk. For forskning kan likevel hele samlingen gjøres tilgjengelig etter søknad og/eller avtale. De språkmodellene – og andre KI-ressurser – som er utviklet av NB, deles fritt. Ikke bare blir språkmodeller brukt i forskning og undervisning; de brukes også av andre organisasjoner både i det offentlige og i privat næringsliv.»¹⁶

Prosjekt Mimir

I desember 2023 fikk Nasjonalbiblioteket i oppdrag å sette i gang et forsknings- og utviklingsprosjekt for å utforske verdien av opphavsrettslig beskyttet materiale i opplæring av norske generative språkmodeller, samt vurdere en eventuell kompensasjonsordning. Relevante norske forskningsmiljøer ble invitert til å delta i prosjektet, og forfatter- og forleggerorganisasjoner ble invitert til å følge prosjektet.

Nasjonalbiblioteket beskriver fasene i prosjektet slik:

- Videreutvikle Nasjonalbibliotekets digitale norske tekstkorpus
- Trene en rekke norske generative språkmodeller på ulike utdrag fra tekstkorpuset
- Evaluere ytelsen til de ulike språkmodellene
- Dokumentere funn og observasjoner.

Nasjonalbiblioteket vil oppsummere sine funn og observasjoner om verdien av beskyttet materiale fra evalueringen av de trente språkmodellene. Ambisjonen er å fullføre prosjektet innen juli 2024.

¹⁶ [Svein Arne Brygfjeld: Kunstig intelligens og Nasjonalbiblioteket \(Stat & Styring, idunn.no\)](#)

De involverte forskningsmiljøene er NorwAI – Norsk senter for AI-innovasjon (NTNU), Institutt for informatikk: Språkteknologigruppen (Universitetet i Oslo) og Sigma2, en organisasjon som administrerer databehandlings- og lagringsressurs for forskere i Norge.

I et møte hos Nasjonalbiblioteket 8. februar ble NBs ståsted, ambisjoner og forståelse av lovverket presentert for interesserte rettighetshavere.

Et viktig spørsmål blir hva man gjør med resultatene av forskningsprosjektene som NB skal gjennomføre våren 2024. Hvis resultatet publiseres i artikkelform, er det åpenbart uproblematisk. Hvis en språkmodell som er skapt i forbindelse med forskning, gjøres tilgjengelig for videre bruk, vil – når digitalmarkedsdirektivet er gjennomført i norsk rett – artikkel 4 være et aktuelt grunnlag, og man kan reservere seg eller inngå en avtale. Om grunnlagsmaterialet for å bygge de enkelte språkmodeller er opphavsrettsbeskyttet eller ikke, er naturligvis viktig med tanke på hva resultatet av forskningen eventuelt kan brukes til senere.

NB opplyser at de modellene ikke er resultatet i forskningsprosjektet, men instrument for å ta fram kunnskap. Det skal altså gjennomføres evaluering av modellene, og det er de resultatene man er ute etter.

Et av sporene til Nasjonalbiblioteket er å utvikle tre ulike språkmodeller: en modell som kun er trent på opphavsrettslig materiale, en «fusion-modell» som kombinerer trening på opphavsrettslig og fritt tilgjengelig materiale og en modell som kun baserer seg på fritt tilgjengelig materiale.

Rettigheter

Nasjonalbiblioteket garanterer at materialet det trenes på i Mimir ikke blir tilgjengelig for andre. Meta vil ikke kunne se materialet, selv om det er deres Llama-modell som brukes. Arbeidet gjøres i lukkede systemer. NB vet ikke om modellene som skal brukes til å forske på norske modeller, Llama (Meta) og Mistral (Mistral AI), kan være skapt med materiale som er opphavsrettsbeskyttet, uten avtale med rettighetshavere.

Svein Arne Brygfjeld, som leder KI-prosjektene på Nasjonalbiblioteket, sier til KI©K-prosjektet:

“For vår del har vi valgt å ikke videreformidle samlinger med data fra Internett uten at vi har kontroll på rettigheter. Men vi har valgt å følge samme praksis som andre som trener modeller, det vil si at slike data brukes for trening.

Et interessant perspektiv er forholdet mellom treningsdata og de trente modellene med hensyn til rettigheter. Etter det jeg forstår, er det etter hvert etablert en generell oppfatning/forståelse innen jussen i Europa at modellene kan distribueres og brukes uavhengig av rettighetsforholdene i treningsdata.”

Nasjonalbiblioteket har laget en oversikt over hva de mener de kan bruke til tradisjonell forskning og til forskning på KI-modeller/språkteknologi:

Kilder	Åpent innhold	Kan brukes til forskning	Kan brukes til språkteknologi
Bøker under opphavsrett	Nei	Ja	Nei
Bøker i det fri	Ja	Ja	Ja
Aviser i det fri	Ja	Ja	Ja
Aviser under opphavsrett, uten avtaler	Nei	Ja	Nei
Aviser med avtaler	Nei	Ja	Delvis
NRK-artikler	Nei	Ja	Nei
VG Debatt	Nei	Delvis	Nei
Twitter	Nei	Delvis	Nei
Reddit	Nei	Delvis	Nei
Facebook	Nei	Delvis	Nei
Wikipedia	Ja	Ja	Ja
Målfrid (statlige nettsider)	Ja	Ja	Ja
Stortinget, offentlige rapporter	Ja	Ja	Ja
Lovdata	Delvis	Ja	Delvis

Kompensasjonsordning

Som en oppfølging av prosjekt Mímir har Kulturdepartementet bedt Nasjonalbiblioteket å vurdere grunnlaget for en eventuell kompensasjonsordning for norske rettighetshavere og eventuelt utvikle et forslag til en slik ordning.¹⁷

Kopinor er bedt om å koordinere «en mindre gruppe», bestående av forfatter- og forleggerorganisasjoner, for å følge prosjektet. Gruppen er nå dannet, og består av representanter for fem forfatterorganisasjoner, Forleggerforeningen og Kopinor. I tillegg har NB på eget initiativ invitert Mediebedriftenes Landsforening inn i en lignende rolle.

Denne gruppen har ingen formell rolle eller tydelig mandat i prosjektet, men Nasjonalbiblioteket virker interessert i å sikre støtte fra rettighetshaverne for deres anbefaling til Kulturdepartementet. Det er også antydnet at Nasjonalbiblioteket primært ønsker å utforske en modell med en kollektiv avtalelisensavtale. De gode erfaringene med Bokhylla-avtalen er en av grunnene til at Nasjonalbiblioteket har fått dette oppdraget av Kulturdepartementet.

Det er positivt at myndighetene med dette ser behovet for å utrede en kompensasjon for bruk av opphavsrettslig beskyttet materiale. Den videre innretningen på en slik eventuell kompensasjonsordning bør være et prioritert arbeidsområde for rettighetshaverne i 2024.

¹⁷ [Skal undersøke bruk av opphavsrettslig beskyttet materiale i trening av norsk kunstig intelligens \(Regjeringen.no\)](#)

Tillegg

Tidslinje

2022	
30. november	OpenAI lanserer ChatGPT
2023	
7. februar	Microsoft integrerer ChatGPT i Bing-søk. På det tidspunkt har allerede mer enn 100 millioner mennesker testet ChatGPT.
17. januar	Getty Images kunngjør at de saksøker Stability AI for å bruke opphavsrettsbeskyttede bilder til å trene sin KI
14. mars	OpenAI lanserer GPT-4. Tre og en halv måned etter lansering av ChatGPT kom OpenAI med en kraftig forbedret versjon, den nye GPT-4. Dette var en multimodal modell som kunne håndtere både bilder og tekst.
21. mars	Google slipper sin egen generativ KI-chatbot, Bard
23. mars	Storbritannia oppdaterer sitt AI Policy white paper som første gang kom i august 2022. I en enda senere versjon heter det blant annet at Storbritannia vil beholde sin status som en global leder innen KI, blant annet relatert til opphavsrettslovgivning og generativ KI, samt sikre riktige balanse mellom å beskytte rettighetshavere og kreative industrier, samtidig som KI-utviklere får adgang til dataene de trenger.
Vinteren og våren	Google og Microsoft begynner å slippe generative KI-funksjoner over sine eksisterende forretningsapplikasjoner. En jevn strøm av regelmessig oppdaterte utgivelser av KI-verktøy i andre medier: for eksempel Canva tekst-til-bilde designverktøy, Midjourney tekst-til-video og ElevenLabs tekst-til-tale
29. mars	Elon Musk og over 1000 teknologiledere og KI-eksperter ber om en pause i KI-utviklingen til dens effekter kan forstås, forutses og kontrolleres bedre
Sommeren	Nyhetsleverandører over hele verden begynner å ta offentlig stilling til generativ KI og etablerer etiske rammeverk for bruk. Publisister i USA og Europa understreker behovet for menneskelig tilsyn med KI-genererte produkter, og mediehus eksperimenterer med å bruke KI-verktøy. Medieorganisasjoner begynner å ta stilling til teknologiplattformer, spesielt rundt rettferdig bruk av deres innhold som treningsdata.
13. juli	Associated Press signerer en avtale med OpenAI for å lisensiere nyhetshistorier.
Høsten	Myndigheter i svært mange land begynner å ta stilling til, eller drøfte, regulatoriske utfordringer. Samtidig starter mediehus med å blokkere for ytterligere trening på deres internettarkiver.
Oktober	Amerikanske Actors Guild varsler massesøksmål mot Open AI og Microsoft for brudd på opphavsretten ved trening av KI-modellene sine. Komiker og forfatter Sarah Silverman anlegger et eget søksmål mot Open AI.
30. oktober	Presidentordre i USA tar til orde for en samarbeidsbasert tilnærming til å jobbe med store teknologiselskaper om sikkerhet og trygghet, i et forsøk på å balansere risiko opp mot innovasjon.

8. desember	Digitalmarkedsdirektivet innlemmes i EØS-avtalen. Direktivet inneholder de mest omfattende samlede endringene i EUs opphavsretsregulering siden 2001. Endringene i åndsverkloven er planlagt lagt fram for Stortinget i 2025.
9. desember	Enighet mellom EU-rådet og parlamentet om den nye AI Act. Den er restriktiv og krever blant annet at KI-systemer må gjennomgå en vurdering før lansering.
13. desember	Axel Springer offentliggjør en avtale med OpenAI om å åpne for bruk av sitt innhold i treningsdata til språkmodeller og AI-tjenester.
27. desember	The New York Times saksøker OpenAI og Microsoft for opphavsrettsbrudd.
Desember	Google Gemini rulles ut til brukere over hele verden. Gemini er en såkalt multimodal KI-modell, trent til å forstå og generere innhold på tvers av flere strukturer inkludert tekst, bilder, lyd og video. Dette gjør tjenesten i stand til å forstå og løse komplekse problemer ved å integrere informasjon fra ulike datakilder.
2024	
Februar	En amerikansk domstol avviser deler av søksmålet komiker og forfatter Sarah Silverman har anlagt mot Open AI for brudd på opphavsretten.
15. februar	OpenAI lanserer Sora, en tekst-til-video-tjeneste som lager inntil ett minutt lange realistiske videofilmer kun basert på tekst-beskrivelser.
13. mars	EU-parlamentet vedtar den endelige teksten til KI-forordningen (AI Act)

Ordliste

AI (Artificial Intelligence)

AI er det internasjonale uttrykket, brukt første gang i 1956. På norsk har KI og kunstig intelligens fått gjennomslag, mens det er fortsatt ganske vanlig å bruke AI som forkortelse. Våre naboland bruker AI. (I Sverige kan man ikke si KI, for «kunstig intelligens» betyr noe helt annet.)

AI Act

AI Act (norsk: KI-forordningen) er EUs regelverk for kunstig intelligens, som etter planen skal endelig vedtas i april 2024. Se nærmere omtale i rapporten.

Digitalmarkedsdirektivet

Digitalmarkedsdirektivet (Digital Single Market Directive, DSM) ble vedtatt av EU i 2019 og er gjennomført i de enkelte medlemsstatene. Se nærmere omtale i rapporten.

Fair use

Fair use er et prinsipp i USAs lovgivning, som betyr at selv om noe er opphavsrettslig beskyttet, kan det likevel brukes uten tillatelse fra innehaveren av opphavsretten i bestemte situasjoner. For å avgjøre om en bestemt bruk gjort av et verk anses å være *fair use* vurderes bruken opp mot fire faktorer:

1. Formålet med og karakteren av bruken, blant annet om bruken er kommersiell eller ikke-kommersiell
2. Arten av det opphavsrettsbeskyttede verket
3. Hvor stor del av det beskyttede verket som er brukt
4. Hvilken virkning bruken av verket har på det potensielle markedet for eller verdien av det beskyttede verket.

I norsk og annen kontinentaleuropeisk rett vil mange slike tilfeller i stedet fanges opp av definerte unntak i loven (f.eks. sitatrett og kopiering til privat bruk).

Generativ KI

Generativ KI betegner kunstig intelligens som er i stand til å generere tekst, bilder eller andre data ved hjelp av generative modeller. Per i dag genereres det som svar på «prompts» eller spørsmål/bestillinger.

KI-crawler

Språkmodellene bygges opp ved at en KI-crawler beveger seg rundt på ulike nettsteder og samler sammen tekster, både språklige regler og kunnskap fra innholdet. Forskjellen på en KI-crawler og det som tidligere ble brukt for å hente ut informasjon fra internett, er at den ikke har bruk for at innholdet på nettsidene er strukturert eller samlet samme sted.

Korpus

Korpus (tekstkorpus) er en samling av tekster som brukes til språklige eller litterære undersøkelser. Det kan være en elektronisk database med tekster som er samlet for å representere et språk eller en språkvariant. Korpuset kan inneholde ulike typer tekster, som bøker, aviser, tidsskriftartikler, og transkripsjoner fra tale, og det kan være både generelt, dekke mange temaer og sjangre, eller det kan være mer spesialisert og fokusert på et bestemt emne eller felt.

Korpuset kan være *annotert*, «tagget», med lingvistisk informasjon, som ordklasser eller syntaktisk struktur, noe som gjør det mulig å utføre detaljerte analyser av språkbruk og mønstre. Tekstkorpus

brukes i utviklingen av språkmodeller og de bidrar til å forbedre teknologier som maskinoversettelse, talegjenkjenning, og automatiserte chatbots/digitale assistenter.

NCC (Norwegian Colossal Corpus)

NCC er Nasjonalbibliotekets språkkorpus som kun inneholder ikke opphavsrettsbeskyttet materiale. Det er sist oppdatert i 2021, men skal oppdateres i 2024.

NCC+ er Nasjonalbibliotekets språkkorpus som også inneholder opphavsrettsbeskyttet materiale. Den oppdateres fortløpende, for eksempel får den daglig nytt innhold fra norske medier.

Prompt

Prompt kan oversettes med «ledetekst» og er en instruksjon, spørsmål eller setning som brukes for å initiere en respons eller handling fra en KI-tjeneste. Systemet analyserer prompten og genererer et svar basert på den gitte bestillingen.

Språkmodell

En språkmodell er et beregningsverktøy som kan forstå og generere menneskelig språk ved å identifisere mønstre og sammenhenger i store mengder tekstdata. Slike modeller er trent på store datasett med tekstdata for å identifisere mønstre, sammenhenger og strukturer i språket. Språkmodeller spiller en avgjørende rolle i utviklingen av teknologier som stemmeaktiverte assistenter, automatiserte chatbots for kundeservice, og andre applikasjoner som krever forståelse eller generering av naturlig språk.

Store språkmodeller (Large Language Models, LLM) som GPT-4 m.fl. har tilegnet seg evnen til språkforståelse og språkgenerering ved å lære statistiske sammenhenger fra tekstdokumenter. Takket være denne «opplæringsprosessen» kan en god LLM brukes til generering av tekster, basert på innholdet i språkdatabasen.

Tekst- og datautvinning

Tekst- og datautvinning (Text and Data Mining, TDM) er kjent fra blant annet språkteknologi gjennom flere tiår. Der det finnes store tekstsamlinger kan man forske på ulike deler av tekstene, som for eksempel hvordan forekomsten av ord og uttrykk endrer seg over tid eller annen type tekstanalyse som er vanskelig å fange opp manuelt. Når det nå bygges store språkmodeller er det et langt mer avansert resultat av tekst- og datautvinning enn det man har sett for seg tidligere.

I digitalmarkedsdirektivet defineres tekst- og datautvinning som «*any automated analytical technique aimed at analysing text and data in digital form in order to generate information which includes but is not limited to patterns, trends and correlations.*»

Tretrinnstesten

Tretrinnstesten er en internasjonalt anerkjent ramme for nasjonal handlefrihet når det gjelder adgangen til å gjøre unntak fra eneretten. Tretrinnstesten innebærer at et unntak må oppfylle tre vilkår. Det skal:

1. Bare gjelde i spesielle og avgrensede tilfeller
2. Ikke skade den normale utnyttelse av verket
3. Ikke på urimelig måte tilsidesette opphavsmannens legitime interesser.

Tretrinnstesten finnes formulert i Bernkonvensjonen art. 9 (2) og er gjentatt blant annet i EUs opphavsrettsdirektiv (2001).