



# Machine learning (ML) til sikring af datakvalitet i de nationale sundhedsregistre

Sundhedsdatastyrelsen (SDS) d. 9. oktober 2024  
E-observatoriet 2024

Mikkel Wermer Steen

# Indhold

- **Baggrund, resultater og erfaringer med ML-initiativer til forbedring af datakvalitet**
- **Planer og fremtidigt perspektiv**
- **Afrunding / spørgsmål**



# Baggrund og formål med ML-datakvalitetsprojekter

# Hvad ser vi i SDS?

## 1. Mere kompleks data:

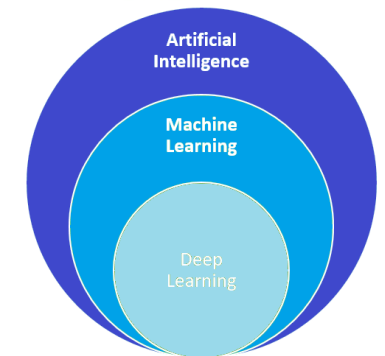
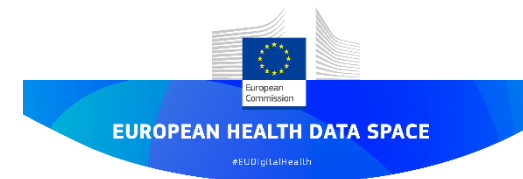
- Vi modtager løbende et større og mere varieret volumen af data
- Datastrukturen bliver stadig mere kompleks

## 2. Øget interesse for sundhedsdata:

- Flere ønsker adgang til data
- Gælder både til primær og sekundær anvendelse

## 3. Nye initiativer på vej:

- EHDS-forordningen er under implementering
- Nationalt set ønskes data mere tilgængeligt
- Fokus er på bedre deling og tilgængelighed af data



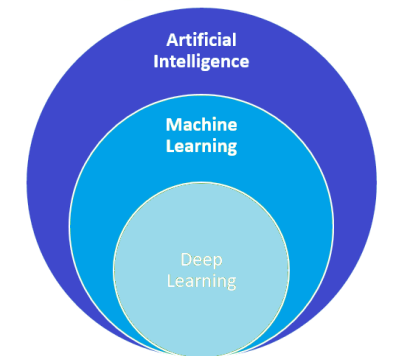
# Hvordan påvirker det SDS?

## 1. Større krav:

- Vi oplever stigende efterspørgsel på oplysninger om data:
  - Tilgængelighed
  - Indhold
  - Datakvalitet

## 2. Behov for ibrugtagen af nye metoder til:

- Videreudvikling og effektivisering af det eksisterende datakvalitetsarbejde
- Kvalitetsarbejde med nuværende og fremtidige centrale sundhedsregistre



# Hvad gør SDS?

## AI-metoder afprøves med henblik på at:

### 1. Identificere fejl og mangler i data

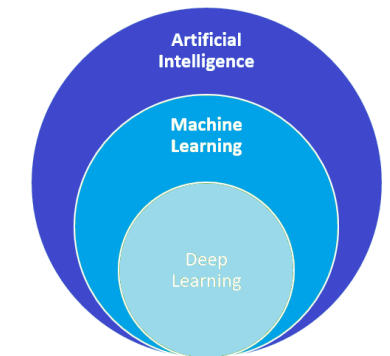
- Se skjulte mønstre i vores sundhedsregistre
- Kunne fange og reagere hvis noget ændrer sig i måden, data indberettes på

### 2. Få overblik over vores data

- Finde ud af i hvilke situationer der typisk forekommer fejl eller mangler

### 3. Opdage problemer tidligt

- Finde fejl, før de bliver et stort problem
- Hurtigt kunne gribe ind med kvalitetsforbedrende tiltag



# Anvendelse af ML til at understøtte datakvalitetsarbejdet

## Succeskriterier for datakvalitetsinitiativerne



### Forbedre datakvalitet

- Identificere ændringer i indberetningsmønstre og indhold
- Afdække fejl og mangler i sundhedsdata
- Få overblik over datakvalitet på tværs af datakilder



### AI kompetenceudvikling

- Opbygge kompetencer internt
- Sidemandsoplæring og gennemsigtighed (undgå "black box")
- Udvikling af generiske metoder som kan benyttes på tværs af flere indberetningsområder

# Datakvalitetsarbejde i SDS

SDS har undersøgt, hvordan ML kan forbedre kvaliteten af sundhedsdata.

**Case 1**  
*Cancer diagnoser*

Validering af  
diagnoseindberetninger til  
Landspatientregistret

*Pilotprojekt afsluttet, plan for videreførelse af  
projektet*

**Case 2**  
*Outlier detection*

Identificering af anomalier i  
indberetningsfrekvenser i  
Landspatientregistret

*Pilotprojekt afsluttet, plan for videreførelse af  
projektet*

**Case 3**  
*Dødsårsager*

Validering af  
dødsårsagskoder i  
Dødsårsagsregistret

*Modeludvikling i afsluttende fase, indledende  
test ift. idriftsættelse er i gang*

Fokus i dag er på validering af dødsårsagskoder (Case 3) i Dødsårsagsregisteret (DAR)



# Case 3 – Dødsårsager

# Proces for bestemmelse af den underliggende dødsårsag



Borger dør



**Dødsattest** Side 1a  
I henhold til afsnit XIII i sundhedsloven, j. LBR nr. 1188 af 24. september 2016 Til pårørende/bestemt

Personnummer (CPR nummer): (Ved dødfald angives moderens CPR-nummer) <input type="text"/>		Mand: <input type="checkbox"/>	Kvinde: <input type="checkbox"/>		
Fulde navn: <input type="text"/>		Sygesikringsgruppe: <input type="text"/>			
Vej: <input type="text"/>		Nr/etage/side: <input type="text"/>			
Postnr: <input type="text"/>	By: <input type="text"/>	Land: <input type="text"/>			
Dødfødt: <input type="checkbox"/>	Fødselsdato: <input type="text"/>	Klokkeslæt: <input type="text"/>	Dreng: <input type="checkbox"/>	Pige: <input type="checkbox"/>	
Dødstidspunkt: Dato: <input type="text"/>		Klokkeslæt: <input type="text"/>			
Findetidspunkt: Dato: <input type="text"/>		Klokkeslæt: <input type="text"/>			
Død på sygehus/hospice (navn): <input type="text"/>					
Afdeling: <input type="text"/>					
Død på kendt adresse (vej): <input type="text"/>		Nr/etage/side: <input type="text"/>			
Postnr: <input type="text"/>	By: <input type="text"/>	Plejehjem: <input type="checkbox"/>	Eget hjem: <input type="checkbox"/>		
Dødssted uden adresse: <input type="text"/>					
Fundet død på kendt adresse (vej): <input type="text"/>		Nr/etage/side: <input type="text"/>			
Postnr: <input type="text"/>	By: <input type="text"/>	Plejehjem: <input type="checkbox"/>	Eget hjem: <input type="checkbox"/>		
Findested uden adresse: <input type="text"/>					
Attestudfyldende læges funktion: Egen læge <input type="checkbox"/> Vagtlæge <input type="checkbox"/> Hospitalslæge <input type="checkbox"/> Overlæge i Styrelsen for Patientsikkerhed <input type="checkbox"/>					
Ligsyn: Dato: <input type="text"/>		Klokkeslæt: <input type="text"/>		Hospicelæge <input type="checkbox"/>	
Dødstegn: Rigor: <input type="checkbox"/>		Liveores: <input type="checkbox"/>	Cadaverositas: <input type="checkbox"/>	Maceratio: <input type="checkbox"/>	Andet: <input type="checkbox"/>
Kontakt til politiet: Ja: <input type="checkbox"/> Nej: <input type="checkbox"/>					
Elektroniske implantater: Ja - og de(t) er fjernet: <input type="checkbox"/> Ja - og de(t) er ikke fjernet: <input type="checkbox"/> Nej: <input type="checkbox"/> Ved ikke: <input type="checkbox"/>					
<small>A. Underrettede læge har syet liget af ovenstående og forefundet ovenstående dødstegn. Bekræfter, at der ikke foreligger tilkrytningsforhold, som omhandlet i Justitsministeriets bekendtgørelse om lægers adgang til at konstatere dødens indtræden, foretage ligsyn, udstede dødsattest og foretage obduktion § 1.</small>		<small>B. Underrettede læge har syet liget af ovenstående og forefundet ovenstående dødstegn. Bekræfter, at der ikke foreligger tilkrytningsforhold, som omhandlet i Justitsministeriets bekendtgørelse om lægers adgang til at konstatere dødens indtræden, foretage ligsyn, udstede dødsattest og foretage obduktion § 1.</small>			
<small>Har ikke fundet ovenstående af den i sundhedsloven § 179, stk. 1 nævnte art, og der er eller vil ske således ikke grund til mistanke om, at døden er forårsaget ved en forbyrdelse, jf. Kirkeministeriets bekendtgørelse om begravelse og ligbrænding § 2 stk. 2.</small>		<small>Har fra politiet modtaget meddelelse om, at offentlig indberetning efter sundhedslovens § 179, stk. 1, ikke giver anledning til retsmedicinsk ligsyn, og at der intet er til hinder for, at ligbrænding kan finde sted, jf. Kirkeministeriets bekendtgørelse om begravelse og ligbrænding § 2 stk. 2.</small>			
Lægens underskrift, navn og adresse, evt. stempel		Lægens underskrift, navn og adresse, evt. stempel			
<small>C. Underrettede overlæge i Styrelsen for Patientsikkerhed (retslæge) har ved retslægeagt ligsyn fundet dødstegn eller andre forhold der er uforenelige med livets betingelser. Bekræfter, at der ikke foreligger tilkrytningsforhold, som omhandlet i Justitsministeriets bekendtgørelse om lægers adgang til at konstatere dødens indtræden, foretage ligsyn, udstede dødsattest og foretage obduktion § 1.</small>		<small>D. Politiets påtegning. Det bekræftes i medfør af sundhedslovens § 182, stk. 2, at der intet er til hinder for at liget begravnes, brændes eller - efter udstedelse af ligsyn - føres ud af landet.</small>			
Overlæge i Styrelsen for Patientsikkerhed (retslægers) underskrift, navn og adresse, evt. stempel		Politets stempel			

Læge udfylder dødsattestens side 1



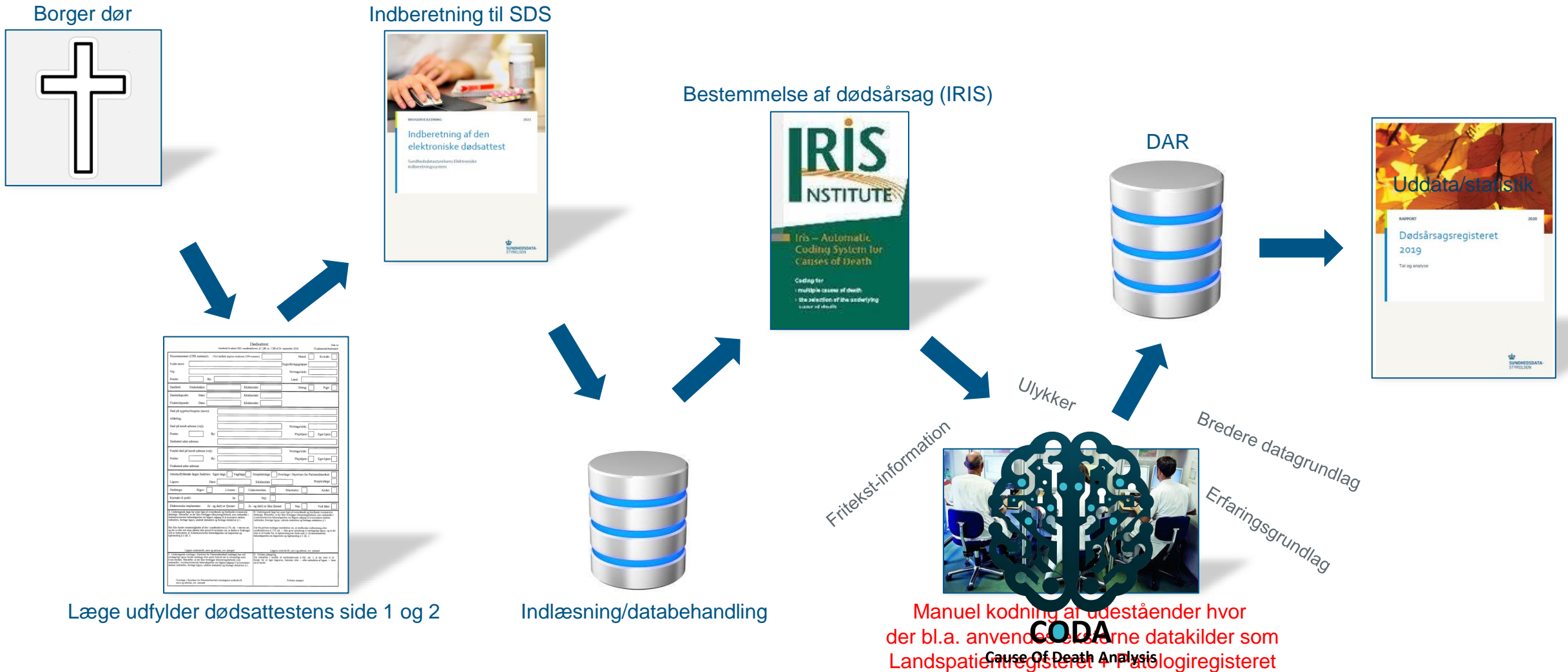
**Dødsattest** Side 2  
I henhold til afsnit XIII i sundhedsloven, j. lov nr. 546 af 24. juni 2005

0395006/2

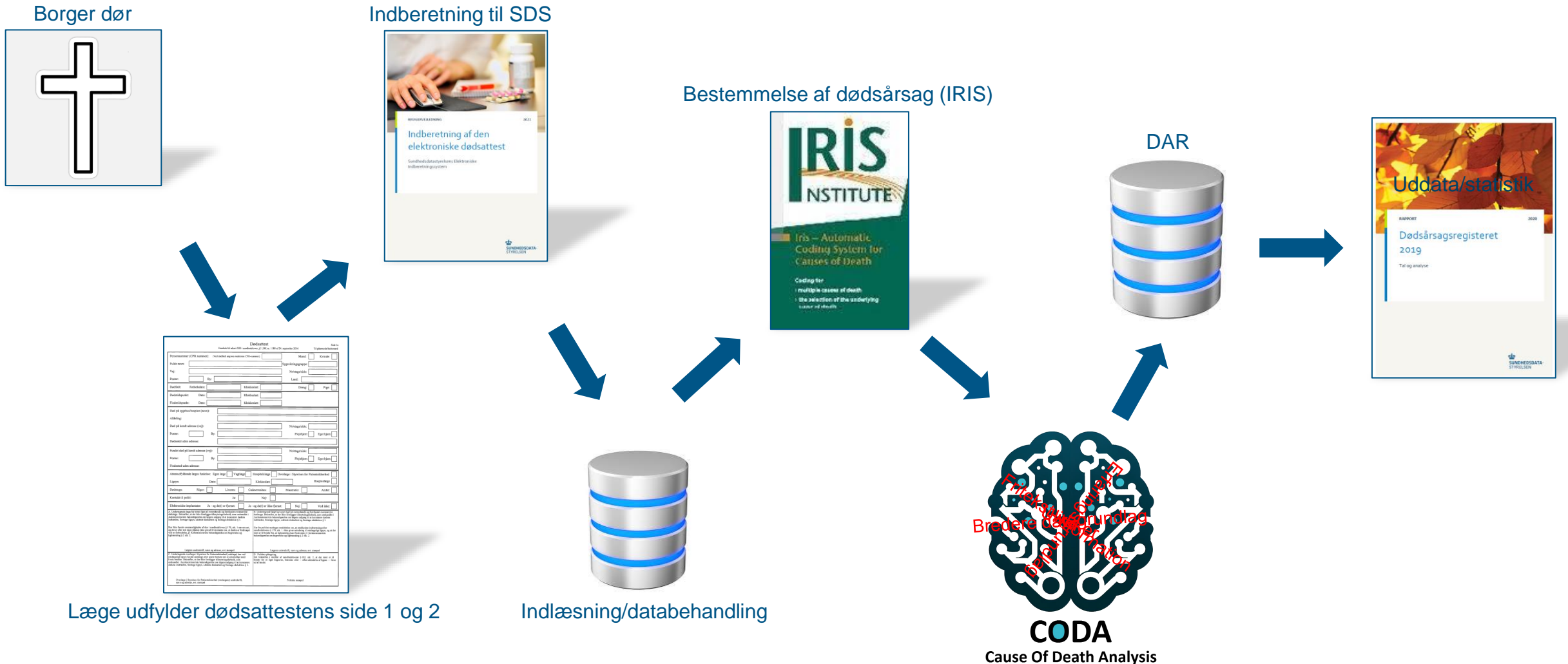
Personnummer (CPR nummer): (Ved dødfald angives moderens CPR-nummer) <input type="text"/>		Mand: <input type="checkbox"/>	Kvinde: <input type="checkbox"/>						
Fulde navn: <input type="text"/>									
Dødsårsag: <input type="checkbox"/> Naturlig død <input type="checkbox"/> Ulykke <input type="checkbox"/> Selvmord <input type="checkbox"/> Drab/Vold <input type="checkbox"/> Uoplyst									
<b>Dødsårsag I.</b> Det sygdoms-, misbrugs- og/eller skadestofbetonede forløb der førte til døden		ICD10	Tidsrum mellem sygdommens opståen og dødens indtræden						
A Den umiddelbare dødsårsag:									
B Som var en følge af:									
C Som var en følge af:									
D Den tilgrundliggende dødsårsag:									
<b>Dødsårsag II.</b> Andre stadig aktive sygdomme, misbrug eller skader, der kan have medvirket til døden		ICD10	Tidsrum mellem sygdommens opståen og dødens indtræden						
Medicin, i forbindelse med forgiftning, medicinbivirkning og misbrug									
Medicinsk præparat (handelsnavn): <input type="text"/>		ATC:							
Medicinsk præparat (handelsnavn): <input type="text"/>		ATC:							
Medicinsk præparat (handelsnavn): <input type="text"/>		ATC:							
Hændelsessted ved ikke-naturlig død:									
Transport-område <input type="checkbox"/>	Bolig-område <input type="checkbox"/>	Produktions-område <input type="checkbox"/>	Handels-område <input type="checkbox"/>	Skole og institution <input type="checkbox"/>	Sports-område <input type="checkbox"/>	Forlystelse og parkområde <input type="checkbox"/>	Fri natur <input type="checkbox"/>	Hav- og søområde <input type="checkbox"/>	Andet <input type="checkbox"/>
Obduktion: Ingen: <input type="checkbox"/> Forbudt: <input type="checkbox"/> Retslig: <input type="checkbox"/> Anden: <input type="checkbox"/>									
Væsentlige obduktionsfund									
<input type="text"/>									
Hvis yderligere undersøgelser er foretaget, specificerer resultaterne heraf									
<input type="text"/>									
Supplerende oplysninger									
<input type="text"/>									
Dato		Attesterende læges navn og stempel							

Læge udfylder dødsattestens side 2

# Proces for bestemmelse af den underliggende dødsårsag



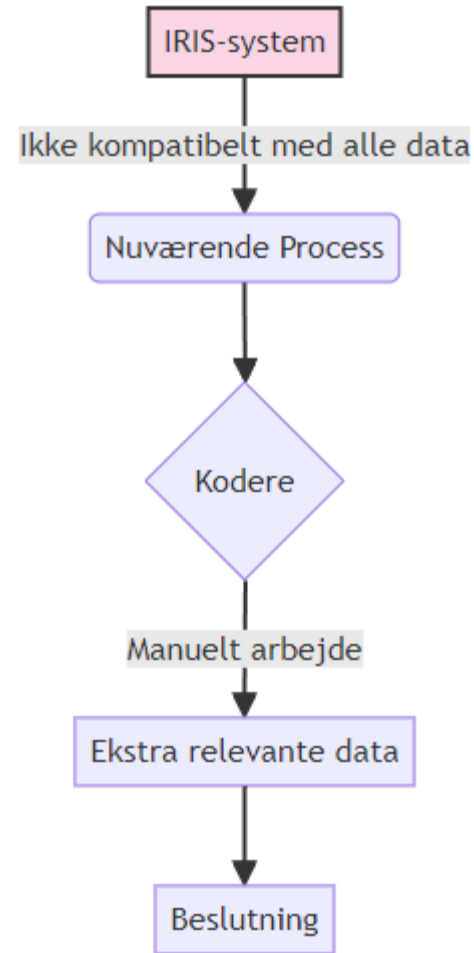
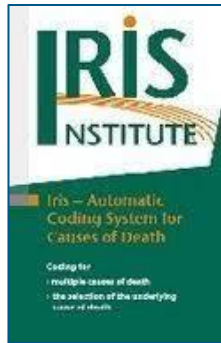
# Proces for bestemmelse af den underliggende dødsårsag



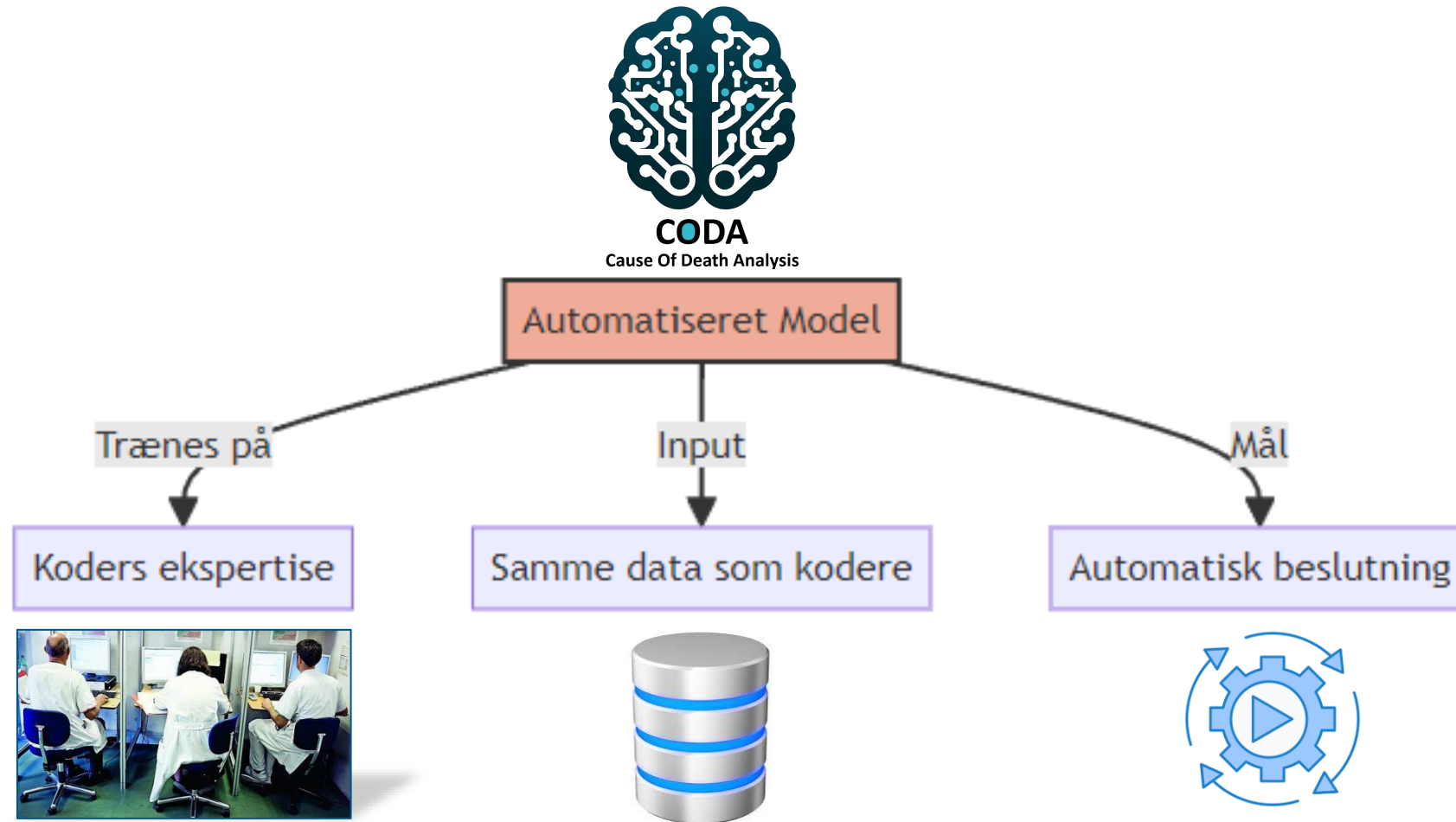
Læge udfylder dødsattestens side 1 og 2

Indlæsning/databehandling

# Proces for bestemmelse af den underliggende dødsårsag



# Proces for bestemmelse af den underliggende dødsårsag

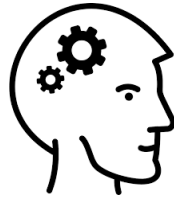


# Proces for bestemmelse af den underliggende dødsårsag



Kodernes centrale rolle

Ekspertter



Træning af model



CODA  
Cause Of Death Analysis

Validering af resultater







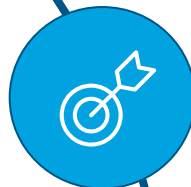
## Problemstilling



### Et datadrevet automatiseret beslutningsredskab til dødsårsager udestår

Det nuværende automatiske kodningsprogram IRIS kræver manuel intervention i 12-18% af indberettede dødsattester for at bestemme den underliggende dødsårsag. Der ønskes udviklet en ML-baseret metode til at effektivisere og understøtte denne manuelle proces

## Mål



Udvikling af en automatiseret model til præcis bestemmelse af den underliggende dødsårsag i tilfælde, hvor IRIS-systemet viser sig utilstrækkeligt eller upålideligt

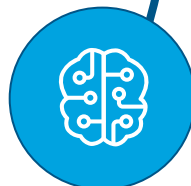
### Datakilder for modellen:

1. Dødsårsagsregistret (DAR):
  - Dødsattestens side 1 og side 2
2. Landspatientregistret (LPR):
  - Patientkontakter, diagnoser, procedurer, forløbselementer, resultatindberetninger
3. Landsregistret for Patologi (PATO):
  - Diagnoser (SNOMED)
4. Oplysninger anvendt i den manuelle kodningsproces
  - Interviews (vægte)

## Data

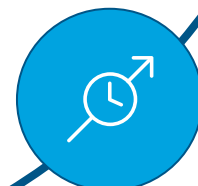


## Analyser



Multiklassifikationsalgoritme der kombinerer registerdata med maskinel analyse af evt. suppleret lægefaglig fritekst om dødsårsagen via en dansk-medicinsk sprog-model\*

## Mål og idéer



Algoritmen trænes på et datasæt bestående af tilfælde, hvor IRIS-systemet har givet usikre resultater. I disse tilfælde er den underliggende dødsårsag blevet manuelt fastlagt af fageksperter



# Matrix på individniveau

## Data til algoritme-træning:

- > 16.000 rækker (dødsattester) i et 80/20 split.
- > 12.000 kolonner med ophav fra: DAR, LPR, PATO og sprogmodel (fritekst)
  - Behandlede kategoriske og numeriske data fra registrene
  - Behandlede vektor-repræsentationer af semantisk forståelse (embeddings) fra sprogmodellen

## Grov eksemplificering af opsætning:

CPR	DAR Side 1 Alder ved død	DAR Side 2 Antal diagnoser	LPR Antal kontakter med hovedspecialet i klinisk onkologi	IRIS-diagnose	Ekspert- diagnose	AI-diagnose	AI-diagnose- sandsynlighed fra modellen
123456-7891	71	10	1	W16	W18	W18	0,90
123456-7892	25	0	12	C619	C619	C618	0,20
123456-7893	76	3	15	D468	D468	D469	0,60
...	...	...	...	...	...	...	...

X\_train/X\_test

y\_train/y\_test

Model estimation

# Matrix på individniveau

## Data til algoritme-træning:

- > 16.000 rækker (dødsattester) i et 80/20 split.
- > 12.000 kolonner med ophav fra: DAR, LPR, PATO og sprogmodel (fritekst)
  - Behandlede kategoriske og numeriske data fra registrene
  - Behandlede vektor-repræsentationer af semantisk forståelse (embeddings) fra sprogmodellen



# Sprogmodel

"Terminal lungekræft med lever- knogle - binyre samt hud-metastaser, i lindrende behandling.  
Indlægges til yderligere lindrende behandling."

"Occult cancer med svær anæmi. Pt. ikke ønsket udredning.  
Sidste tid med pneumoni, så en konkurrerende dødsårsag."

"Var erklæret helbredt. Fik epileptisk anfald som var forårsaget af metastaser fra mamma cancer. Kemoterapi uden

Sprogmodel (Bidirectional Encoder Representations from Transformers - BERT)

Kontekstuelle vektor-repræsentationer af tekstuel information vedr. dødsårsagen

Repræsentation af tekstuel ensartethed

Input til modellen

Får modellen til at "forstå" og anvende tekst-informationen på samme måde som kodeekspertoerne

# XGBoost

## Træ-baseret model ("beslutningstræ")

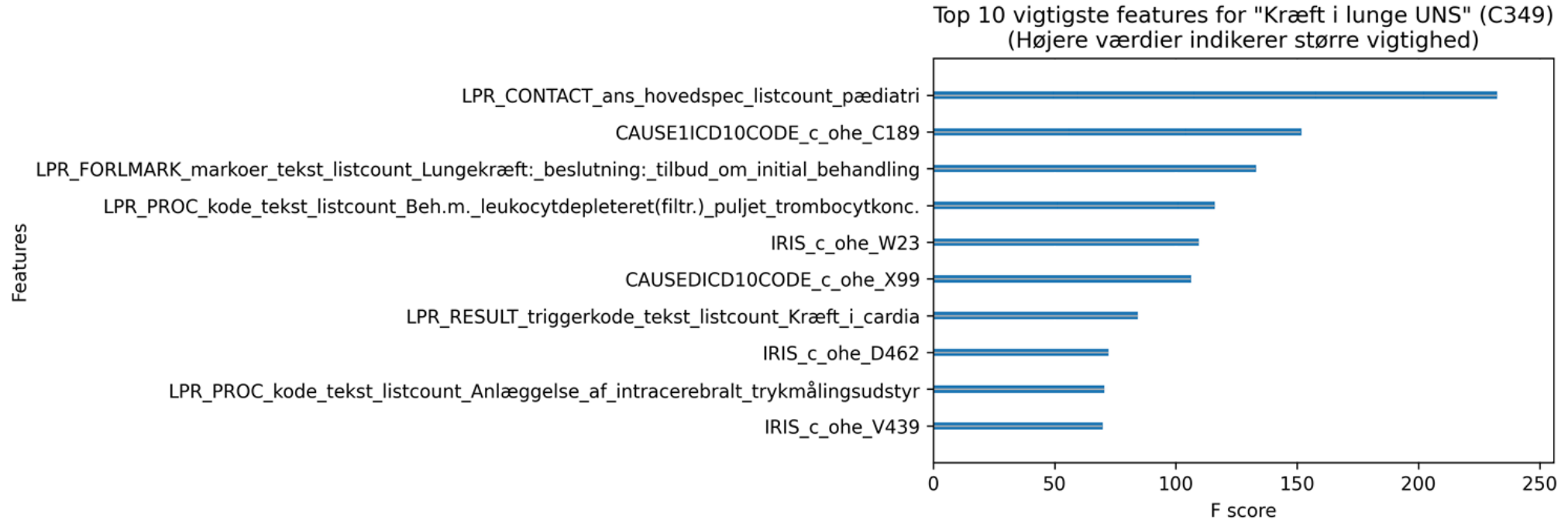
- Vigtigt at kunne gennemgående forstå og forklare modellens beslutninger
- En høj grad af gennemsigtighed gør det lettere at identificere og rette fejl

### Udsnit af den samlede træstruktur for "Kræft i lunge UNS" (C349)

Mange træer bidrager til den endelige beslutning i en kompleks algoritme



## Modellens feature importance ("forklaringselementer")

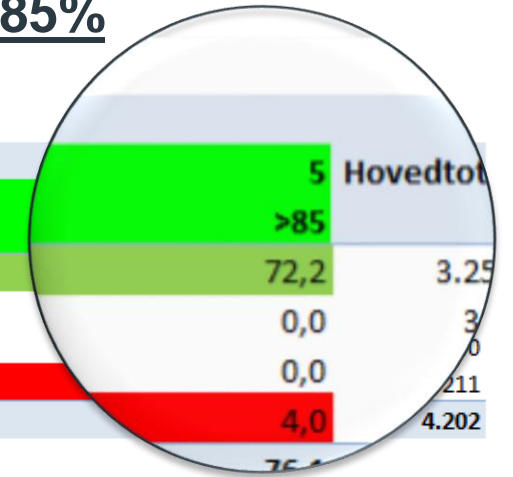


# Modelperformance

- For 72,2% af 4202 attester i testdata var modellens angivne sandsynlighed på >85%

## ICD10 klassifikation - som % af hovedtotal

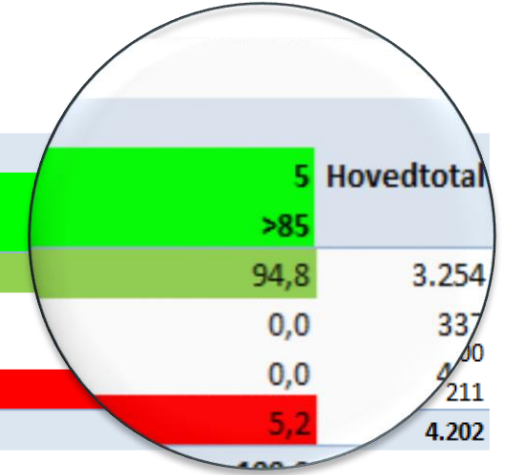
Sum af freq		ai_sandsynlighed_prc_intv_no		ai_sandsynlighed_prc_intv_txt				Hovedtotal
		1	2	3	4	5		
AI_haandtering_ICD10	ai_match	<25	25-50	50-75	75-85	>85		
AI_handling	1	0,0	0,0	0,0	5,3	72,2		3.25
Man_kontrl_ej_match		2,3	3,1	2,6	0,0	0,0		3,0
Man_kontrl_match	1	0,8	2,5	6,2	0,0	0,0		9,5
Model_kontrl	0	0,0	0,0	0,0	1,0	0,0		1,0
Hovedtotal		3,1	5,6	8,9	6,3	4,0		4.202



- Af dødsattesterne med høj angivet modelsandsynlighed, var 94,8% af modellens bud på dødsårsagen korrekt klassificeret ift. den manuelle håndtering

## ICD10 klassifikation - som % af gruppetotal

Sum af freq		ai_sandsynlighed_prc_intv_no		ai_sandsynlighed_prc_intv_txt				Hovedtotal
		1	2	3	4	5		
AI_haandtering_ICD10	ai_match	<25	25-50	50-75	75-85	>85		
AI_handling	1	0,0	0,0	0,0	83,5	94,8		3.254
Man_kontrl_ej_match		74,8	55,1	29,6	0,0	0,0		337,5
Man_kontrl_match	1	25,2	44,9	70,4	0,0	0,0		140,5
Model_kontrl	0	0,0	0,0	0,0	16,5	0,0		16,5
Hovedtotal		100,0	100,0	100,0	100,0	5,2		4.202



# Modelperformance

- For 72,2% af 4202 attester i testdata var modellens angivne sandsynlighed på >85%

## ICD10 klassifikation - som % af hovedtotal

Sum af freq		ai_sandsynlighed_prc_intv_no		ai_sandsynlighed_prc_intv_txt				
		1	2	3	4	5	Hovedtotal	
AI_haandtering_ICD10	ai_match	<25	25-50	50-75	75-85	>85		
AI_handling	1	0,0	0,0	0,0	5,3	72,2		3.254
Man_kontrl_ej_match		2,3	3,1	2,6	0,0	0,0		337
Man_kontrl_match	1	0,8	2,5	6,2	0,0	0,0		400
Model_kontrl	0	0,0	0,0	0,0	1,0	4,0		211
Hovedtotal		3,1	5,6	8,9	6,3	76,1		4.202

- Af dødsattesterne med høj angivet modelsandsynlighed, var 94,8% af modellens bud på dødsårsagen korrekt klassificeret ift. den manuelle håndtering

## ICD10 klassifikation - som % af gruppetotal

Sum af freq		ai_sandsynlighed_prc_intv_no		ai_sandsynlighed_prc_intv_txt				
		1	2	3	4	5	Hovedtotal	
AI_haandtering_ICD10	ai_match	<25	25-50	50-75	75-85	>85		
AI_handling	1	0,0	0,0	0,0	83,5	94,8		3.254
Man_kontrl_ej_match		74,8	55,1	29,6	0,0	0,0		337
Man_kontrl_match	1	25,2	44,9	70,4	0,0	0,0		400
Model_kontrl	0	0,0	0,0	0,0	16,5	5,2		211
Hovedtotal		100,0	100,0	100,0	100,0	100,0		4.202

# Usikkerhed i manuel kodning og modelforbedring

## Afholdte valideringsøvelser (blindtest) viser betydelig variation i manuel kodning

- Usikkerheden på 5% ved høj modelsandsynlighed (>85%), er acceptabel sammenlignet med den menneskelige usikkerhed der er ved manuel kodning af den gruppe af dødsattester som de 5% repræsenterer (som er de særligt svære sager)

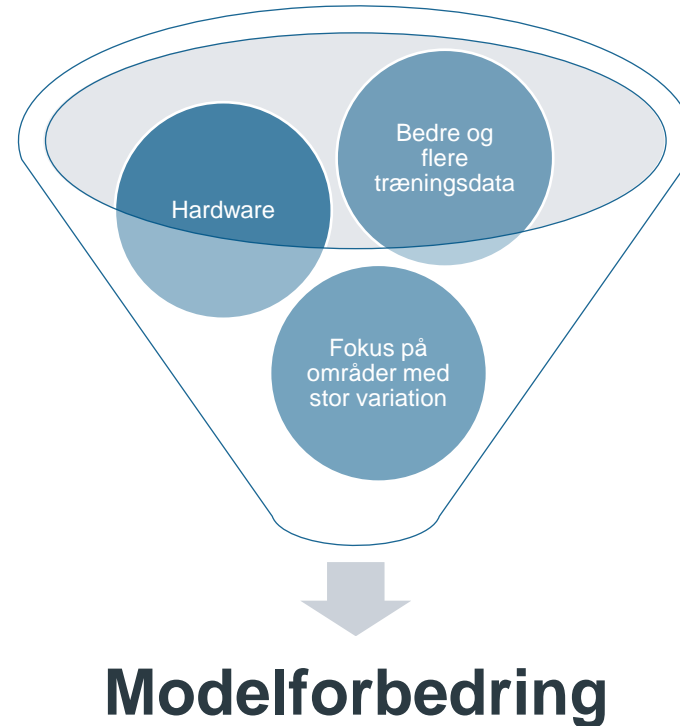




# Variation i manuel kodning og modelforbedring

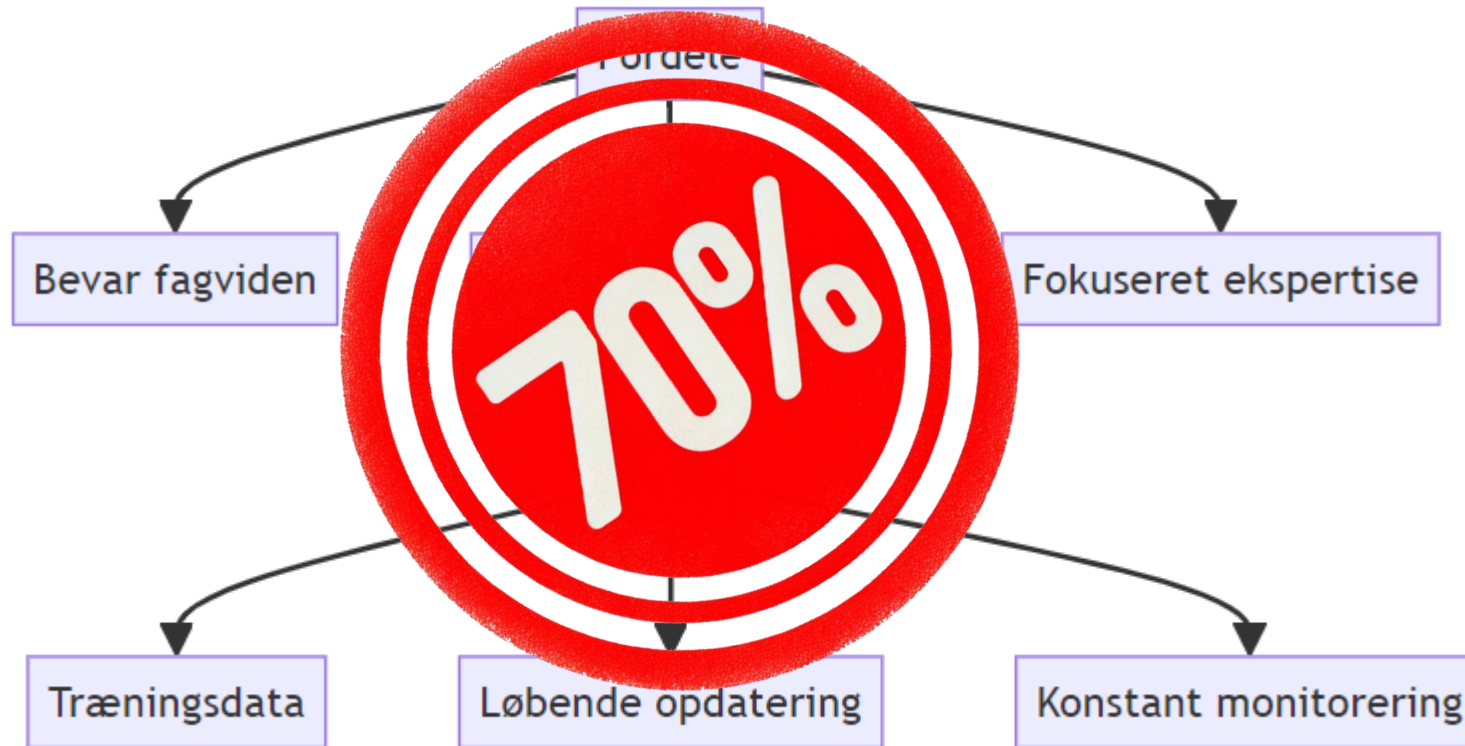
## Afholdte valideringsøvelser (blindtest) viser betydelig variation i manuel kodning

- De 5% forkert klassificerede dødsattester ved høj modelsandsynlighed (>85%), er acceptabelt sammenlignet med variationen i den manuelle kodning for den gruppe af dødsattester som de 5% repræsenterer (som er de særligt svære sager hvor koderne er usikre på dødsårsagen – det vil modellen også være). Modellen er, som koderne, usikker i klassifikationen af de komplicerede dødsattester.



# Model kan fungere som en specialiseret medarbejder

Koderne var ikke umiddelbart imponerede over ML-modellen!



Bare der ikke er flere valideringsøvelser

*Det her giver jo sig selv, modellen kigger jo bare på...*

# Udfordringer og observationer

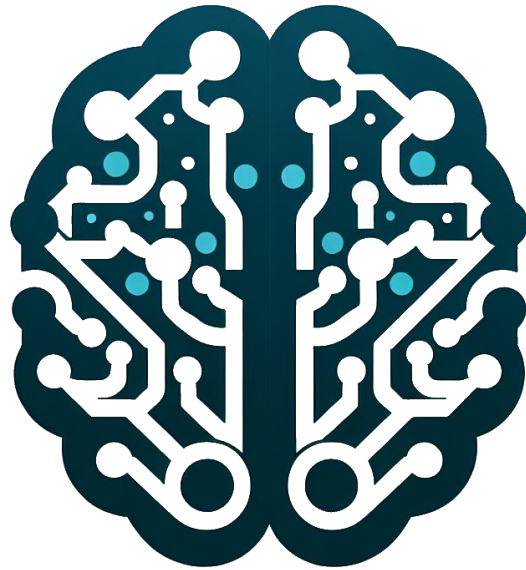
- **Betydelig variation i manuel kodning, især ved øget kompleksitet**
- **Variation i manuel kodning påvirker modellens klassificeringsevne**
- **Lavfrekvente klasser i træningsdata påvirker modellens præstation negativt (håndteres)**
- **Begrænset forbedringspotentiale på modelsiden**
- **Fokus bør være på kvalitet/variation i ekspertvurderingerne (kvalitetsarbejde)**
- **Modelresultater kan udpege områder med størst forbedringspotentiale (målrettet kvalitetsarbejde)**
- **Modelresultater kan påvirke den manuelle kodning (håndteres)**



# Hvor er SDS nu?

## Undersøge mulighederne for anvendelse af modellen som beslutningsværktøj ved høje sandsynligheder

- Målrette kodeeksperternes indsats på at forbedre datakvaliteten – håndtering af særligt komplicerede dødsattester
- Mere målrettet og systematisk arbejde med integrationen af ML-modellen i arbejdet dødsårsagsregisteret forventes at forbedre datakvaliteten i registeret



**CODA**

**Cause Of Death Analysis**



# Opsummering



## Algoritmen kan bestemme den tilgrundliggende dødsårsag

- For **72,2%** af 4202 attester i testdata var modellens angivne sandsynlighed af sit valg på **> 85%**
- Af disse, var **95%** af modellens bud på dødsårsagen korrekt klassificeret ift. den manuelle håndtering
- Resultaterne er kontrolleret ved omfattende valideringsøvelser med kodeeksperterne
- Usikkerheden på **5%** ved høj modelsandsynlighed, er acceptabel sammenlignet med den usikkerhed der er ved manuel kodning af den gruppe af dødsattester som de 5% repræsenterer (de særligt svære sager)



## Algoritmen finder relevante features

- Flere forskellige machine-learning algoritmer er blevet trænet og testet
- XGBoost-algoritmen blev fundet bedst i denne case med en **nøjagtighed<sup>#</sup> på 0,85**
- Algoritmen finder relevante features (explainable AI)

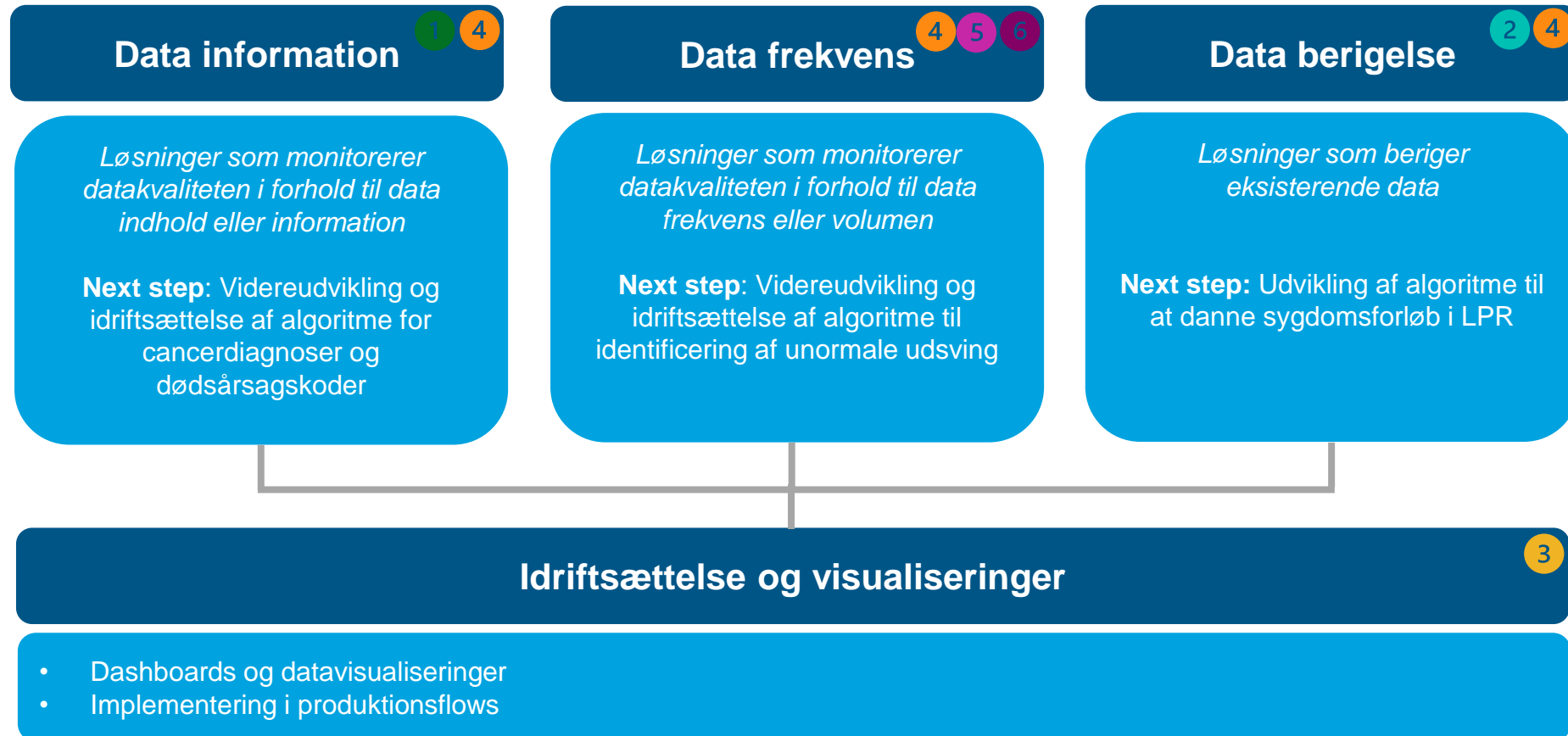


## Algoritmen skal ibrugtages

- Formålet med at undersøge om det er muligt at bestemme den tilgrundliggende dødsårsag via en ML-model er påvist
- Algoritmen vil blive idriftsat som understøttelse til den manuelle ekspertkodning
- *Den ML-baserede løsning har potentiale til at blive et vigtigt redskab i det fremadrettede kvalitetsarbejde for dødsårsagsregisteret*

<sup>#</sup>Nøjagtighed er et udtryk for hvor nøjagtig eller pålidelig algoritmen er. Algoritmen er bedst, jo nærmere nøjagtigheden er på 1.

# Idékatalog med kommende AI-baserede datakvalitetsprojekter i SDS

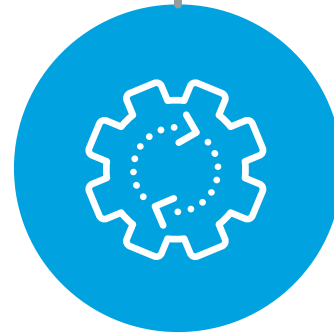


# Foreløbige resultater og perspektiver for videre anvendelse

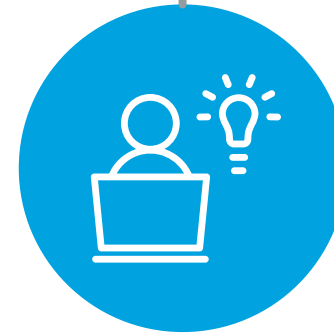
AI har et betydeligt potentiale i SDS og vil kunne anvendes på tværs af organisationen i kvalitetsarbejdet med de centrale sundhedsregistre



AI-metoder kan identificere væsentlige fejl og mangler i indberetninger - også for indberetningsområder hvor der er høj datakvalitet



AI kan automatisere tidskrævende og komplicerede datakvalitetsopgaver



Det er muligt at opbygge kompetencer for at anvende AI-baserede metoder til dataarbejdet i styrelsen



Det kræver målrettet arbejde og ressourcer at opbygge AI-kompetencer og tilhørende infrastruktur



