



SINTEF

Optimering av maskinl ring for  kt b rekraft

Hvem er jeg?

- Adela Nedisan Videsjorden 
- Bachelor i Spr kteknologi (UiO)
- Mastergrad i Programmering og systemarkitektur (UiO)
- I dag: Master of Science (SINTEF – Trustworthy Green IoT)

Dagens tema:

- Veien til b rekraftig maskinl ring



UiO : Universitetet i Oslo



SINTEF

Industri 4.0

- Økning i antall **digitale løsninger** innenfor ulike industriområder.
- Økning i bruk av **sensorer** og **kunstig intelligens (KI)** teknologi.



SINTEF

Industri 4.0

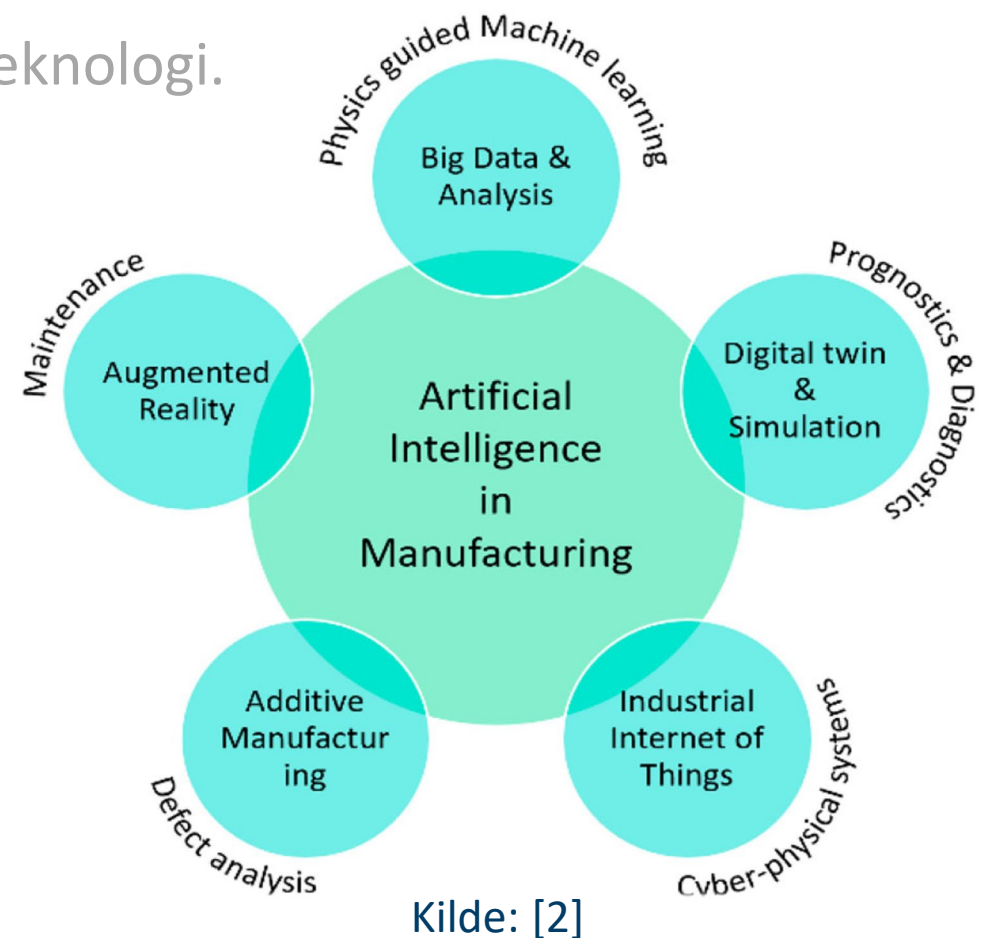
- Økning i antall **digitale løsninger** innenfor ulike industriområder.
- Økning i bruk av **sensorer** og **kunstig intelligens (KI)** teknologi.
- **Manufacturing** (produksjon) er spesielt preget:

“The area of greatest impact by I4.0 is manufacturing”

Kilde: [1]

KI er en hovedingrediens som muliggjør I4.0!

KI-bruk innen produksjon:



Kilde: [2]



SINTEF

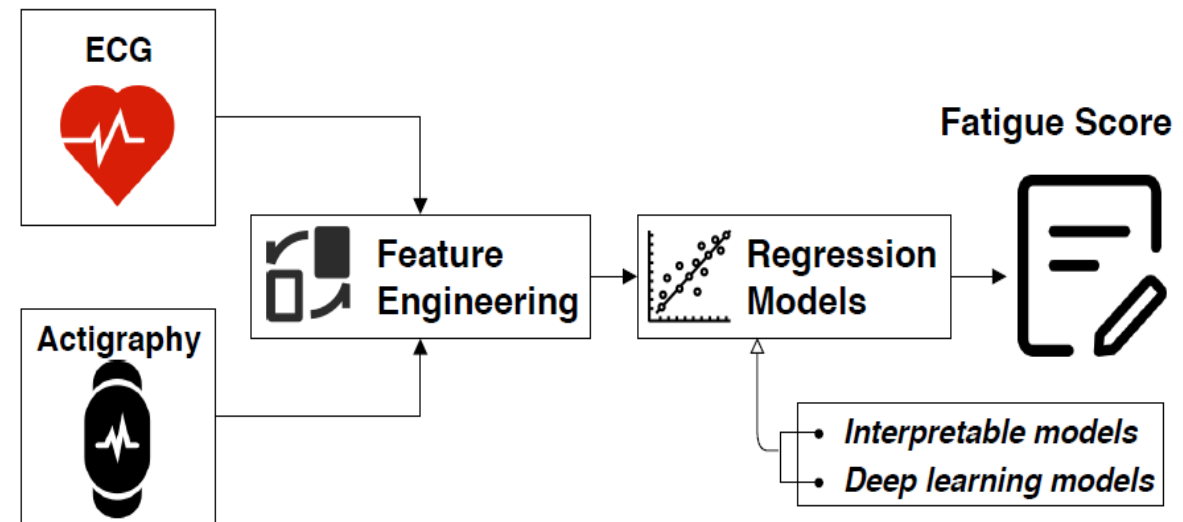
Industri 4.0

- Økning i antall **digitale løsninger** innenfor ulike industriområder.
- Økning i bruk av **sensorer** og **kunstig intelligens (KI)** teknologi.
- **Manufacturing** (produksjon) er spesielt preget.
- Bruk av KI innen **helsesektoren**:

Industri 4.0

- Økning i antall **digitale løsninger** innenfor ulike industriområder.
- Økning i bruk av **sensorer** og **kunstig intelligens (KI)** teknologi.
- **Manufacturing** (produksjon) er spesielt preget.
- **Bruk av KI innen helsesektoren:**

- Støtter spesialister i beslutningsprosessen [3]
- Legemiddelutvikling [4]
- Deteksjon av tretthet, stress og depresjon [5,6,7]



Kilde: [5]



SINTEF

Industri 4.0

- Økning i antall **digitale løsninger** innenfor ulike industriområder.
- Økning i bruk av **sensorer** og **kunstig intelligens (KI)** teknologi.
- **Manufacturing** (produksjon) er spesielt preget.
- Bruk av KI innen **helsesektoren** og **logistikk**:



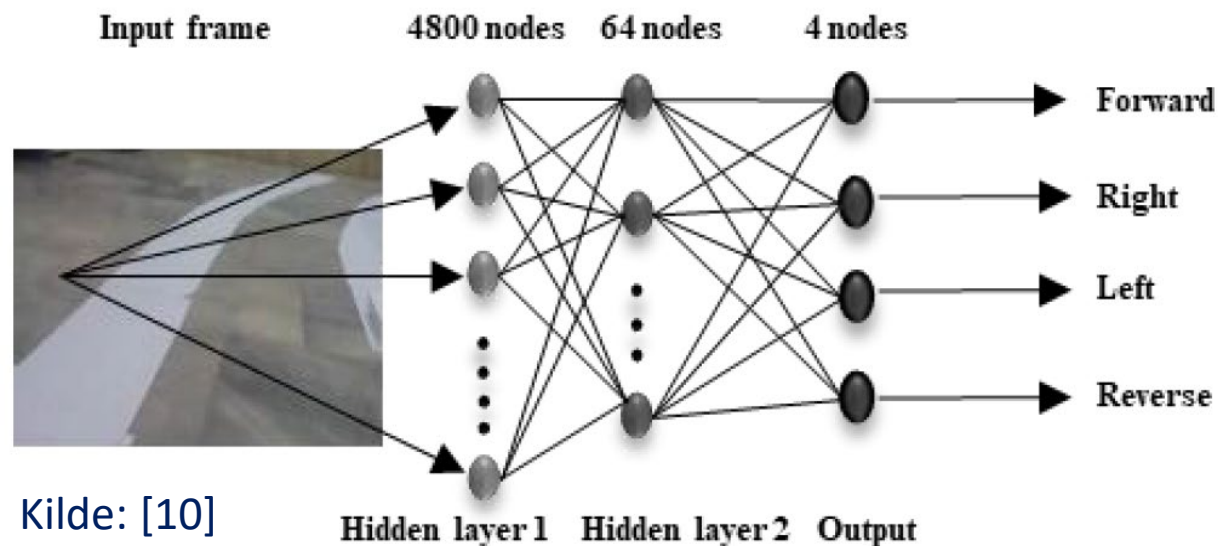
SINTEF

Industri 4.0

- Økning i antall **digitale løsninger** innenfor ulike industriområder.
- Økning i bruk av **sensorer** og **kunstig intelligens (KI)** teknologi.
- **Manufacturing** (produksjon) er spesielt preget.
- Bruk av KI innen **helsesektoren** og **logistikk**:

Ifølge [8], KI har følgende fordeler:

- ✓ **Raskere leveringshastighet**
- ✓ **Høyere pålitelighet**
- ✓ **Lavere driftskostnader**
- ✓ **Forbedret driftseffektivitet**



- Eksempler: pakkelevering via droner [9], selvkjørende biler [10]



SINTEF

Industri 4.0

- Økning i antall **digitale løsninger** innenfor ulike industriområder.
- Økning i bruk av **sensorer** og **kunstig intelligens (KI)** teknologier.
- **Manufacturing** (produksjon) er spesielt preget.
- Bruk av KI innen **helsesektoren** og **logistikk**.

- **Og så videre ...**



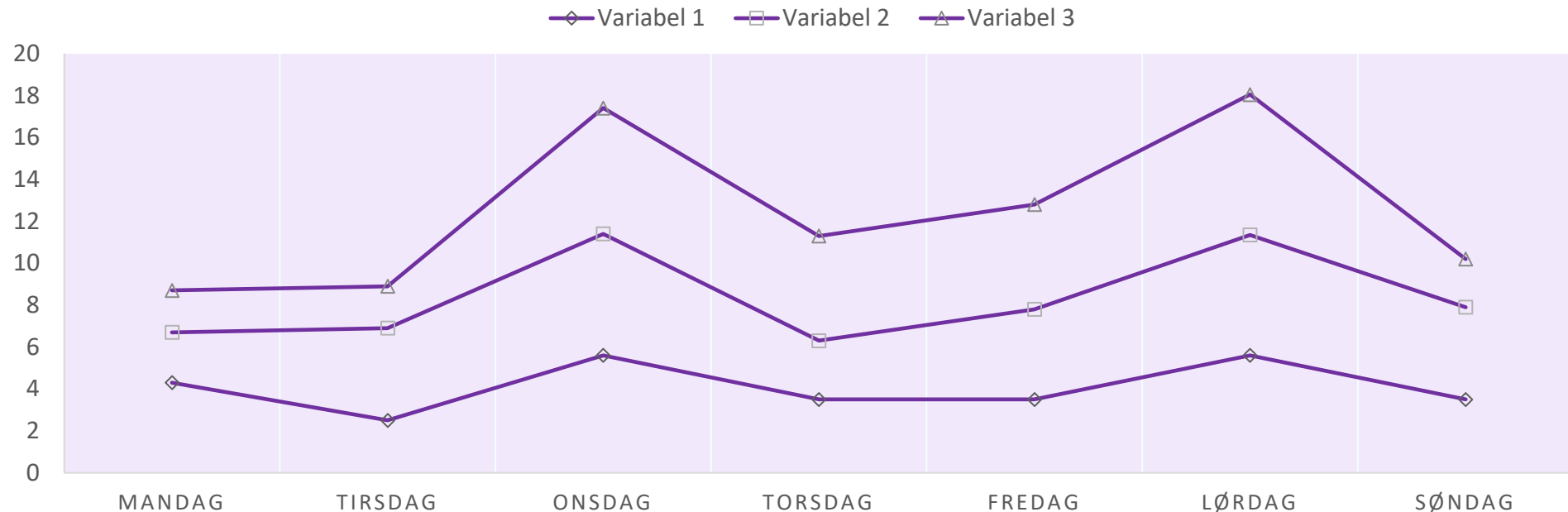
Tidsseriedata (eng: Time-series)

- En **sensor** måler variabler ved bestemte tidsintervaller

Tidsseriedata (eng: Time-series)

- En **sensor** måler variabler ved bestemte tidsintervaller → tidsseriedata

EKSEMPEL

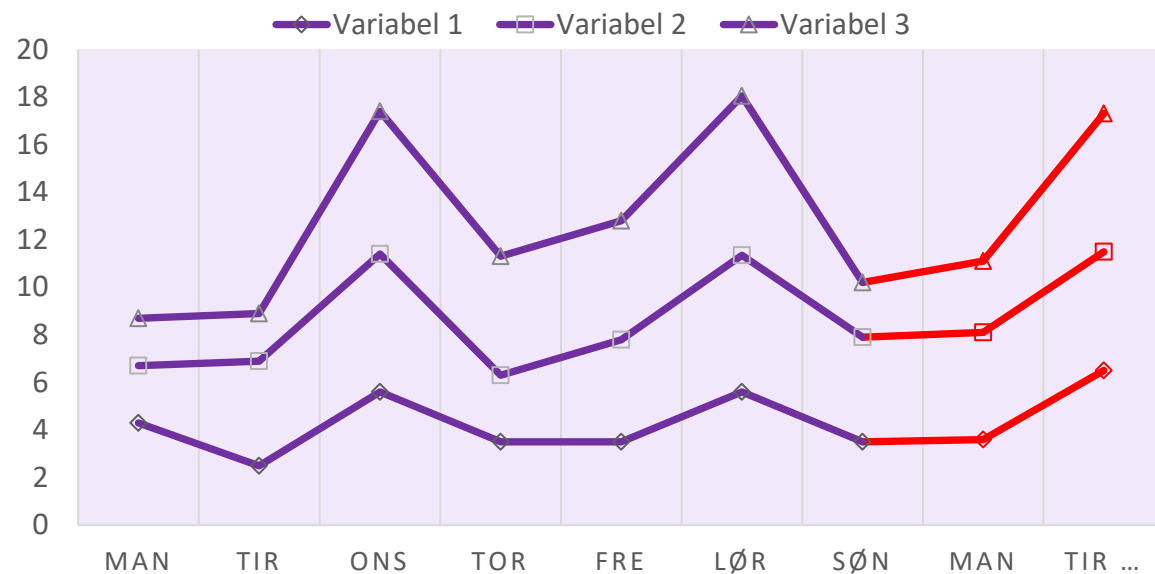


Illustrasjon: sensor som måler 3 ulike variabler en gang om dagen over en periode på en uke.

Tidsseriedata (eng: Time-series)

- En **sensor** måler variabler ved bestemte tidsintervaller → tidsseriedata
- **Maskinlæringsmodeller (ML)** trent på disse dataene kan brukes til å **predikere fremtidige hendelser**

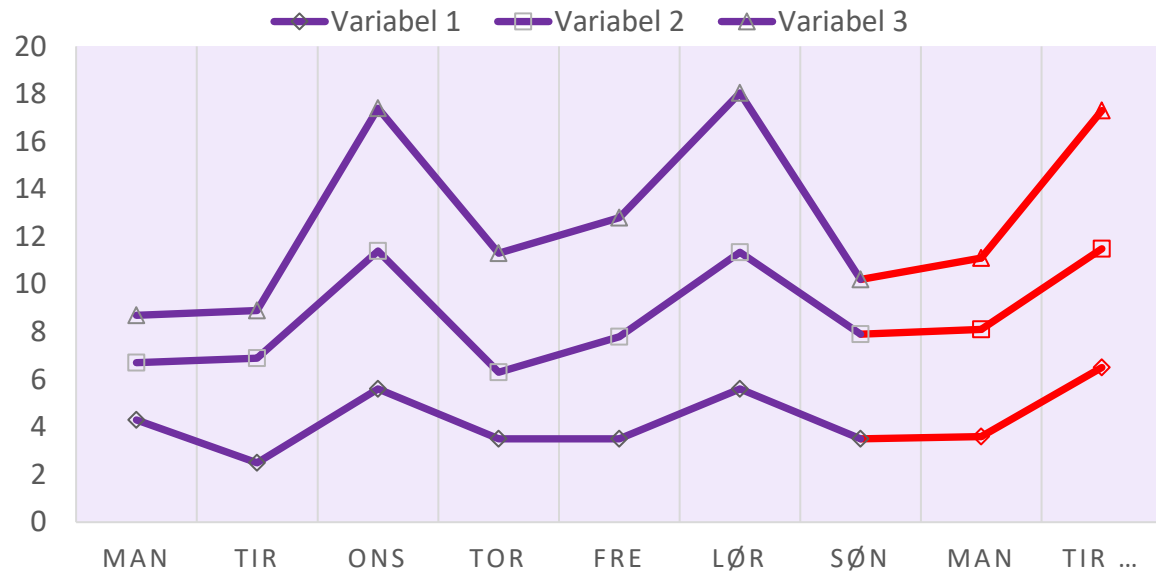
REGRESJON



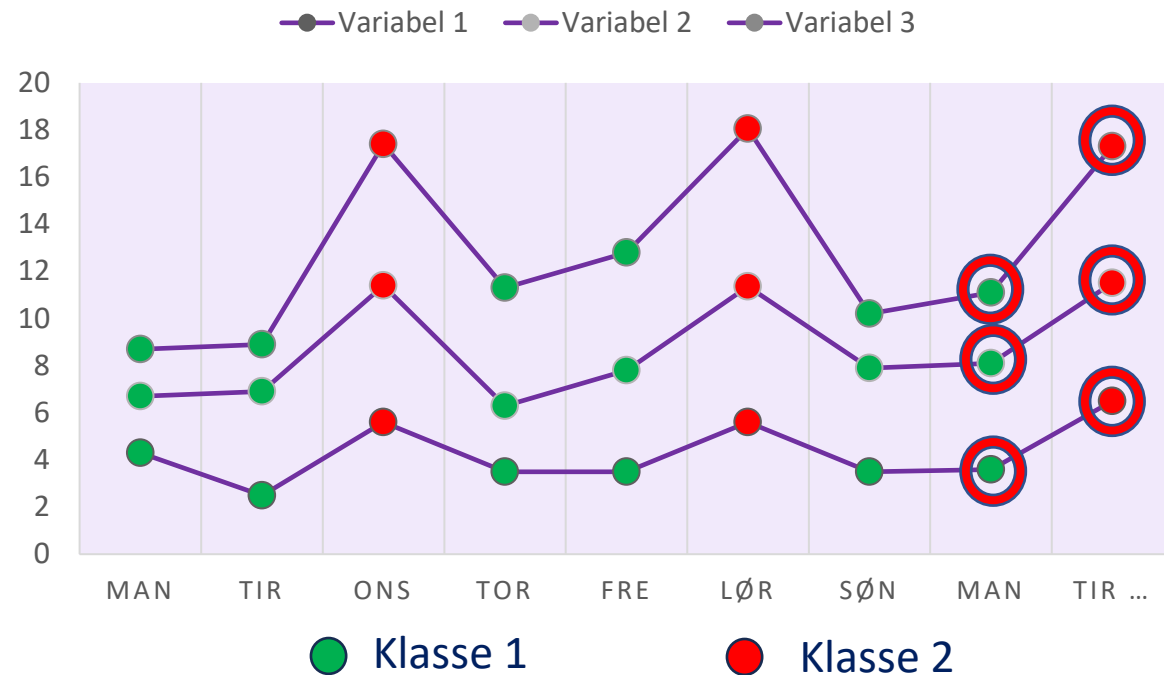
Tidsseriedata (eng: Time-series)

- En **sensor** måler variabler ved bestemte tidsintervaller → tidsseriedata
- **Maskinlæringsmodeller** (ML) trent på disse dataene kan brukes til å **predikere fremtidige hendelser** eller **klassifisere hendelser**.

REGRESJON



KLASSIFIKASJON





Dyplæring (eng: Deep Learning)

- Populært valg: ikke nødvendig å manuell finne viktige trekk i data, sterk ikke-lineær læringsevne, osv.

Dyplæring (eng: Deep Learning)

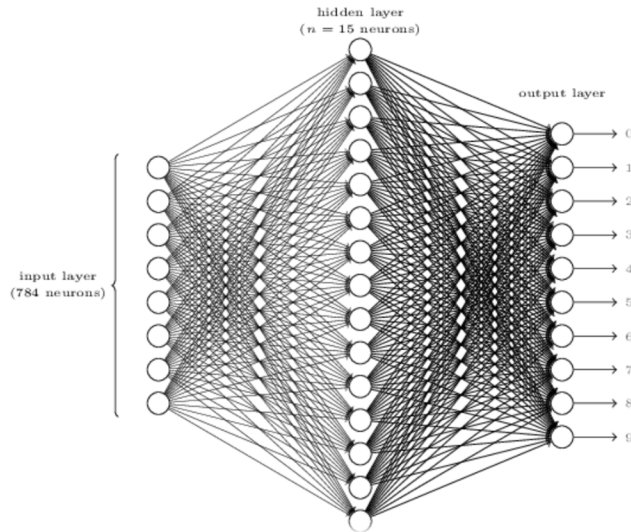
- Populært valg: ikke nødvendig å manuell finne viktige trekk i data, sterk ikke-lineær læringsevne, osv. Flere nettverk **arkitekturer** som kan behandle tidsseriedata!



Dyplæring (eng: Deep Learning)

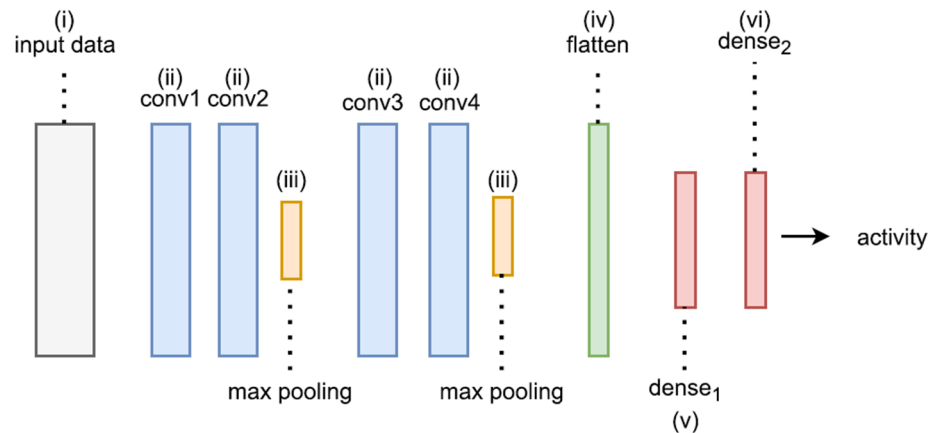
- Populært valg: ikke nødvendig å manuell finne viktige trekk i data, sterk ikke-lineær læringsevne, osv. Flere nettverk **arkitekturer** som kan behandle tidsseriedata!

Multi-Layer Perceptron (MLP)



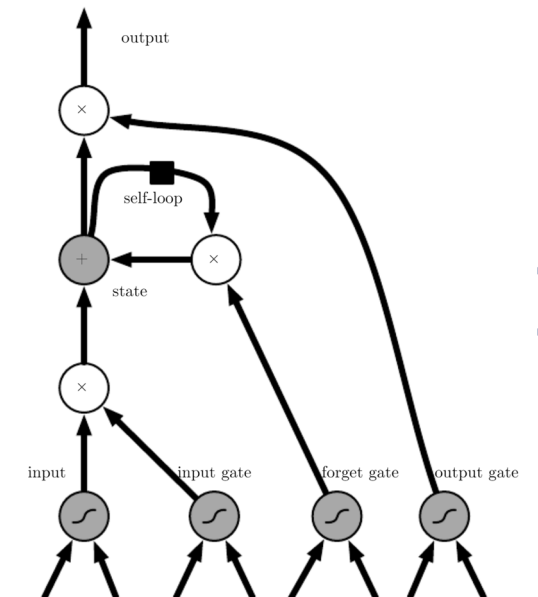
Kilde: [11]

1D Convolutional Neural Network (CNN)



Kilde: [12]

Long Short-Term Memory (LSTM) – celle:



Kilde: [13]



Dyplæring (eng: Deep Learning)

- Populært valg: ikke nødvendig å manuell finne viktige trekk i data, sterk ikke-lineær læringsevne, osv. Flere nettverk **arkitekturer** som kan behandle tidsseriedata!
Skytjenester (eng: Cloud computing) muliggjør **trening** og **inferens** av dyp læring-modeller ved å tilby **ressurselastisitet** [14].



Dyplæring (eng: Deep Learning)

- Populært valg: ikke nødvendig å manuell finne viktige trekk i data, sterk ikke-lineær læringsevne, osv. Flere nettverk **arkitekturer** som kan behandle tidsseriedata!
- **Skytjenester** (eng: Cloud computing) muliggjør **trening** og **inferens** av dyp læring-modeller ved å tilby **ressurselastisitet** [14]. **MEN!** Dette er ikke alltid gjennomførbart pga begrenset nettverksbåndbredde eller av personvernshensyn [15,16].

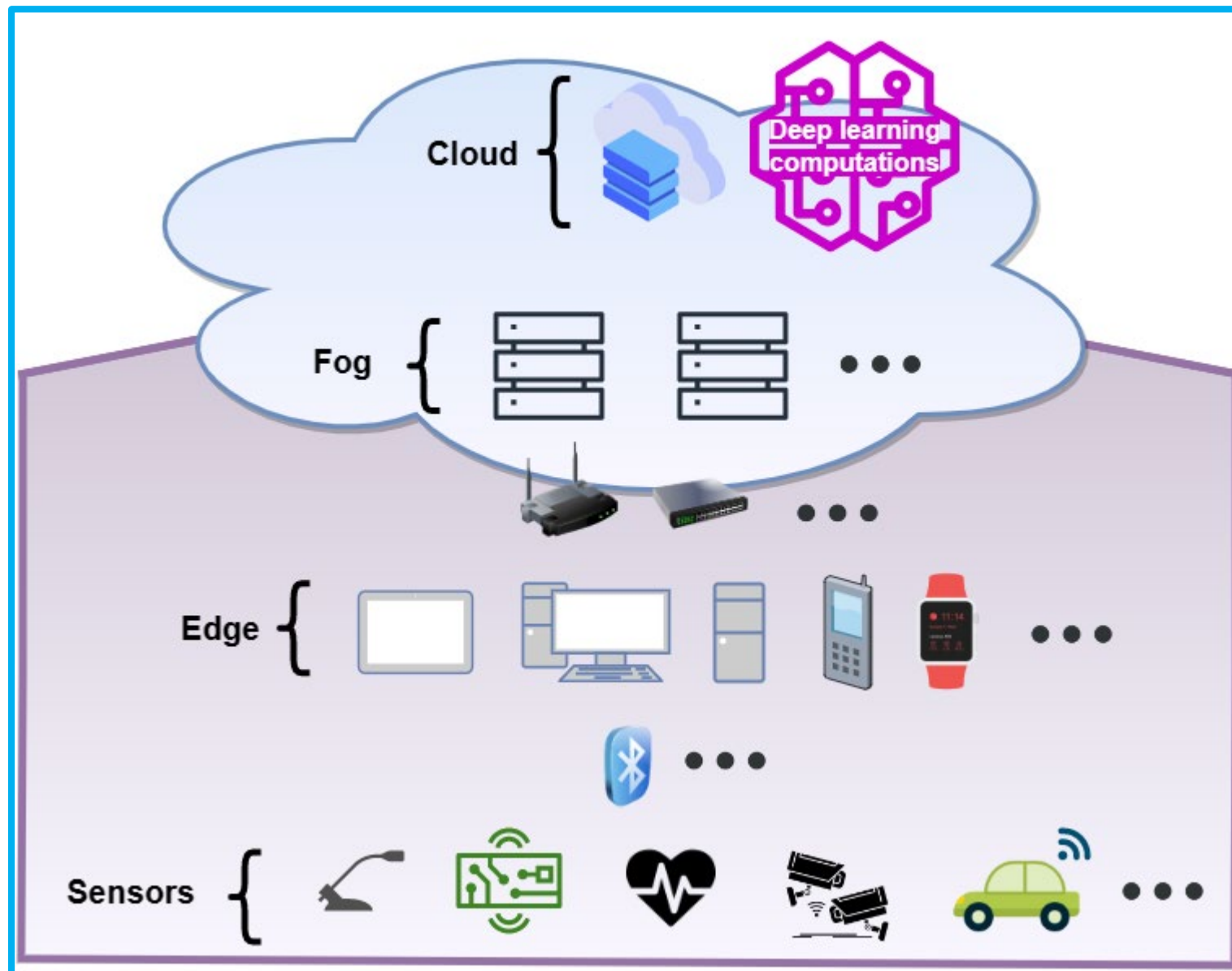
Dyplæring (eng: Deep Learning)

- Populært valg: ikke nødvendig å manuell finne viktige trekk i data, sterk ikke-lineær læringsevne, osv. Flere nettverk **arkitekturer** som kan behandle tidsseriedata!
Skytjenester (eng: Cloud computing) muliggjør **trening** og **inferens** av dyp læring-modeller ved å tilby **ressurselastisitet** [14]. **MEN!** Dette er ikke alltid gjennomførbart pga begrenset nettverksbåndbredde eller av personvernshensyn [15,16].
- **Løsning:** flytt tunge beregninger nærmere **edgen!** [16]



SINTEF

Fra skytjeneste til edge



Dyplæring (eng: Deep Learning)

- Populært valg: ikke nødvendig å manuell finne viktige trekk i data, sterk ikke-lineær læringsevne, osv. Flere nettverk **arkitekturer** som kan behandle tidsseriedata!
Skytjenester (eng: Cloud computing) muliggjør **trening** og **inferens** av dyp læring-modeller ved å tilby **ressurselastisitet** [14]. **MEN!** Dette er ikke alltid gjennomførbart pga begrenset nettverksbåndbredde eller av personvernshensyn [15,16].
- **Løsning:** flytt tunge beregninger nærmere **edgen!** [16]
 - ✓ Behandler data nærmere kilden
 - ✓ Reduserer miljøpåvirkningen knyttet til store datasentre

Også
BÆREKRAFTIG!



Dyplæring (eng: Deep Learning)

- Populært valg: ikke nødvendig å manuell finne viktige trekk i data, sterk ikke-lineær læringsevne, osv. Flere nettverk **arkitekturer** som kan behandle tidsseriedata!
Skytjenester (eng: Cloud computing) muliggjør **trening** og **inferens** av dyp læring-modeller ved å tilby **ressurselastisitet** [14]. **MEN!** Dette er ikke alltid gjennomførbart pga begrenset nettverksbåndbredde eller av personvernshensyn [15,16].
- **Løsning:** flytt tunge beregninger nærmere **edgen!** [16]
- **MEN!** Utfordringer på grunn av begrensede enhetsressurser: beregnings- og minnekrav, strøm- og energiforbruk.



Dyplæring (eng: Deep Learning)

- Populært valg: ikke nødvendig å manuell finne viktige trekk i data, sterk ikke-lineær læringsevne, osv. Flere nettverk **arkitekturer** som kan behandle tidsseriedata!
Skytjenester (eng: Cloud computing) muliggjør **trening** og **inferens** av dyp læring-modeller ved å tilby **ressurselastisitet** [14]. **MEN!** Dette er ikke alltid gjennomførbart pga begrenset nettverksbåndbredde eller av personvernshensyn [15,16].
- **Løsning:** flytt tunge beregninger nærmere **edgen!** [16]
- **MEN!** Utfordringer på grunn av begrensede enhetsressurser: beregnings- og minnekrav, strøm- og energiforbruk.
- **Optimaliseringer** blir derfor anvendt både på **maskinvare-** og **programvarenivå** [17].



SINTEF

Dyplæring (eng: Deep Learning)

- Populært valg: ikke nødvendig å manuell finne viktige trekk i data, sterk ikke-lineær læringsevne, osv. Flere nettværk arkitekturer som kan handle tidsseriedata!
- Skytjenester (eng: Cloud computing) muliggjør trening og inferens av dyp læring-modeller ved å tilby ressursene som ikke alltid gjennomførbart pga begrenset nettverksåndbredde og maskinsyn [15,16].
- **Løsning:** flytt tunge beregninger nærmere data! [16]
- **MEN!** Utfordringer på grunn av ressurskrav: beregnings- og minnekrav, strøm- og energiforbruk.
- **Optimaliseringer** blir derfor anvendt både på maskinvare- og programvarenivå [17].

«En-passer-alle»

optimalisering

på tvers av nettverk

arkitektur?



Modell optimisering

- **Hvorfor?**
- Eliminere unødvendig optimaliseringsinnsats for modeller som håndterer sensor data.



Modell optimisering

- **Hvorfor?**
- Eliminere unødvendig optimaliseringsinnsats for modeller som håndterer sensor data.
- **Hvordan?**
- Finn minimalt sett med **optimaliseringsteknikker** som oppfyller følgende to kriterier:
 - ✓ Redusere ressursutnyttelsen.
 - ✓ Bevare nøyaktigheten.



Modell optimisering

- **Hvorfor?**
- Eliminere unødvendig optimaliseringsinnsats for modeller som håndterer sensor data.
- **Hvordan?**
- Finn minimalt sett med **optimaliseringsteknikker** som oppfyller følgende to kriterier:
 - ✓ Redusere ressursutnyttelsen.
 - ✓ Bevare nøyaktigheten.
- Vi bruker **tidligere forskningsbidrag** til å designe et **rammeverk**. Vi søker etter potensielle **optimaliseringer**, men også **arkitekturer, enheter** og måter å **måle ressursbruk**.

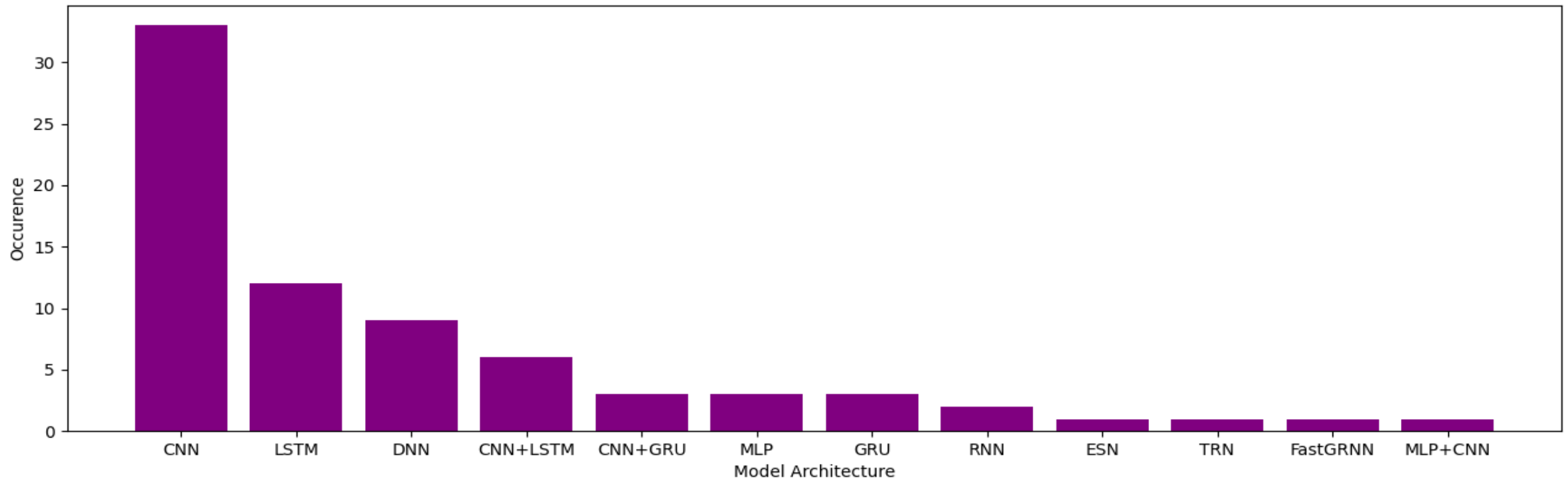


Kort om funn rundt tidsserieanalyse

- **Formål:** helseapplikasjoner utgjør det mest utbredte bruksområdet.

Kort om funn rundt tidsserieanalyse

- **Formål:** helseapplikasjoner utgjør det mest utbredte bruksområdet.
- **Arkitektur:** de mest populære modellarkiteturene er CNN, LSTM, MLP.



Kort om funn rundt tidsserieanalyse

- **Formål:** helseapplikasjoner utgjør det mest utbredte bruksområdet.
- **Arkitektur:** de mest populære modellarkiteturene er CNN, LSTM, MLP.
- **Enhet:** mange eksperimenter kjører på Raspberry Pi eller lignende enheter.



Kort om funn rundt tidsserieanalyse

- **Formål:** helseapplikasjoner utgjør det mest utbredte bruksområdet.
- **Arkitektur:** de mest populære modellarkiturene er CNN, LSTM, MLP.
- **Enhet:** mange eksperimenter kjører på Raspberry Pi eller lignende enheter.
- **Målinger:** det finnes felles indikatorer for ressursbruk/prediksjonsytelse.

ID	Name	Description
KPI₁	Precision / Macro average precision	Correctness of the proportion of identifications in a model
KPI₂	Recall / Macro average recall	Correctness of the proportion of actual positives being correctly identified by the model
KPI₃	Harmonic Precision-Recall Mean (F1) / Macro average F1	Combines precision and recall
KPI₄	Inference Latency	Measure speed during one complete inference request
KPI₅	Parameters	Estimates memory requirements and model complexity
KPI₆	Compressed file size	Estimates necessary memory requirements for a compressed model.
KPI₇	Non-volatile memory	Portion of space allocated for non-volatile storage of model
KPI₈	CPU usage	Estimates system-wide CPU utilization as a percentage
KPI₉	Power consumption	Quantifies the power needed to perform a complete inference request
KPI₁₀	Energy consumption	Quantifies the energy needed to perform a complete inference request

Nøyaktighet (under KPI₁-KPI₃)

Ressursbruk (under KPI₄-KPI₁₀)



Kort om funn rundt tidsserieanalyse

- **Formål:** helseapplikasjoner utgjør det mest utbredte bruksområdet.
- **Arkitektur:** de mest populære modellarkiturene er CNN, LSTM, MLP.
- **Enhet:** mange eksperimenter kjører på Raspberry Pi eller lignende enheter.
- **Målinger:** det finnes felles indikatorer for ressursbruk/prediksjonsytelse.
- **Optimiseringer:** kan fordeles i tre kategorier:
- 1) Maskinvare og distribusjons optimiseringer.

Hardware and distribution optimizations

- Parallelism
- Hardware accelerators
- Memory access optimization
- Near-threshold technology/computing
- Specialized hardware
- Distributed computing



Kort om funn rundt tidsserieanalyse

- **Formål:** helseapplikasjoner utgjør det mest utbredte bruksområdet.
- **Arkitektur:** de mest populære modellarkiturene er CNN, LSTM, MLP.
- **Enhet:** mange eksperimenter kjører på Raspberry Pi eller lignende enheter.
- **Målinger:** det finnes felles indikatorer for ressursbruk/prediksjonsytelse.
- **Optimiseringer:** kan fordeles i tre kategorier:
 - 1) Maskinvare og distribusjons optimiseringer.
 - 2) Optimisering før og mens modellen trener.

Training-time optimizations

- **Keep models small by reducing input size, architecture search.**
- **Architectural design choices: binary models, implementing specific architectures, using other types of operations, attention modules, bottleneck blocks.**
- **Training-time quantization**



Kort om funn rundt tidsserieanalyse

- **Formål:** helseapplikasjoner utgjør det mest utbredte bruksområdet.
- **Arkitektur:** de mest populære modellarkiturene er CNN, LSTM, MLP.
- **Enhet:** mange eksperimenter kjører på Raspberry Pi eller lignende enheter.
- **Målinger:** det finnes felles indikatorer for ressursbruk/prediksjonsytelse.
- **Optimiseringer:** kan fordeles i tre kategorier:
 - 1) Maskinvare og distribusjons optimiseringer.
 - 2) Optimisering før og mens modellen trener.
 - 3) Optimisering etter at modellen er trent.

Inference-time optimizations

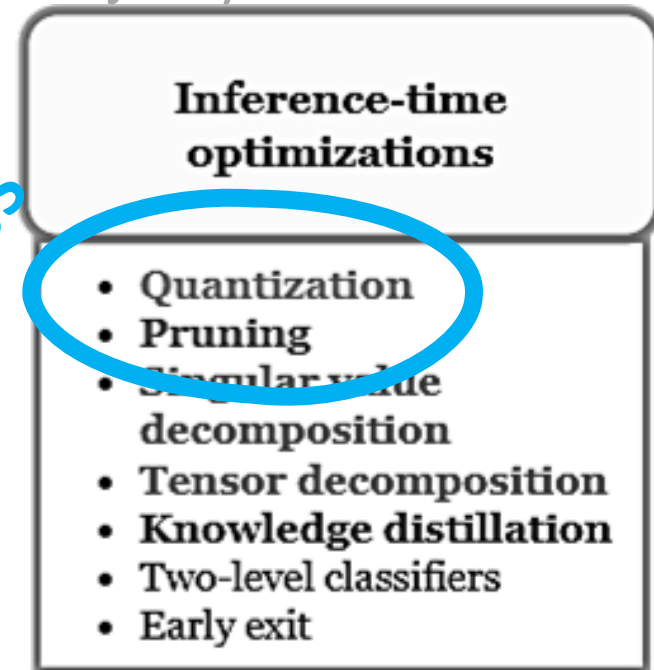
- Quantization
- Pruning
- Singular value decomposition
- Tensor decomposition
- Knowledge distillation
- Two-level classifiers
- Early exit



Kort om funn rundt tidsserieanalyse

- **Formål:** helseapplikasjoner utgjør det mest utbredte bruksområdet.
- **Arkitektur:** de mest populære modellarkiturene er CNN, LSTM, MLP.
- **Enhet:** mange eksperimenter kjører på Raspberry Pi eller lignende enheter.
- **Målinger:** det finnes felles indikatorer for ressursbruk/prediksjonsytelse.
- **Optimiseringer:** kan fordeles i tre kategorier:
 - 1) Maskinvare og distribusjons optimiseringer.
 - 2) Optimisering før og mens modellen trener.
 - 3) Optimisering etter at modellen er trent.

FOKUS





SINTEF

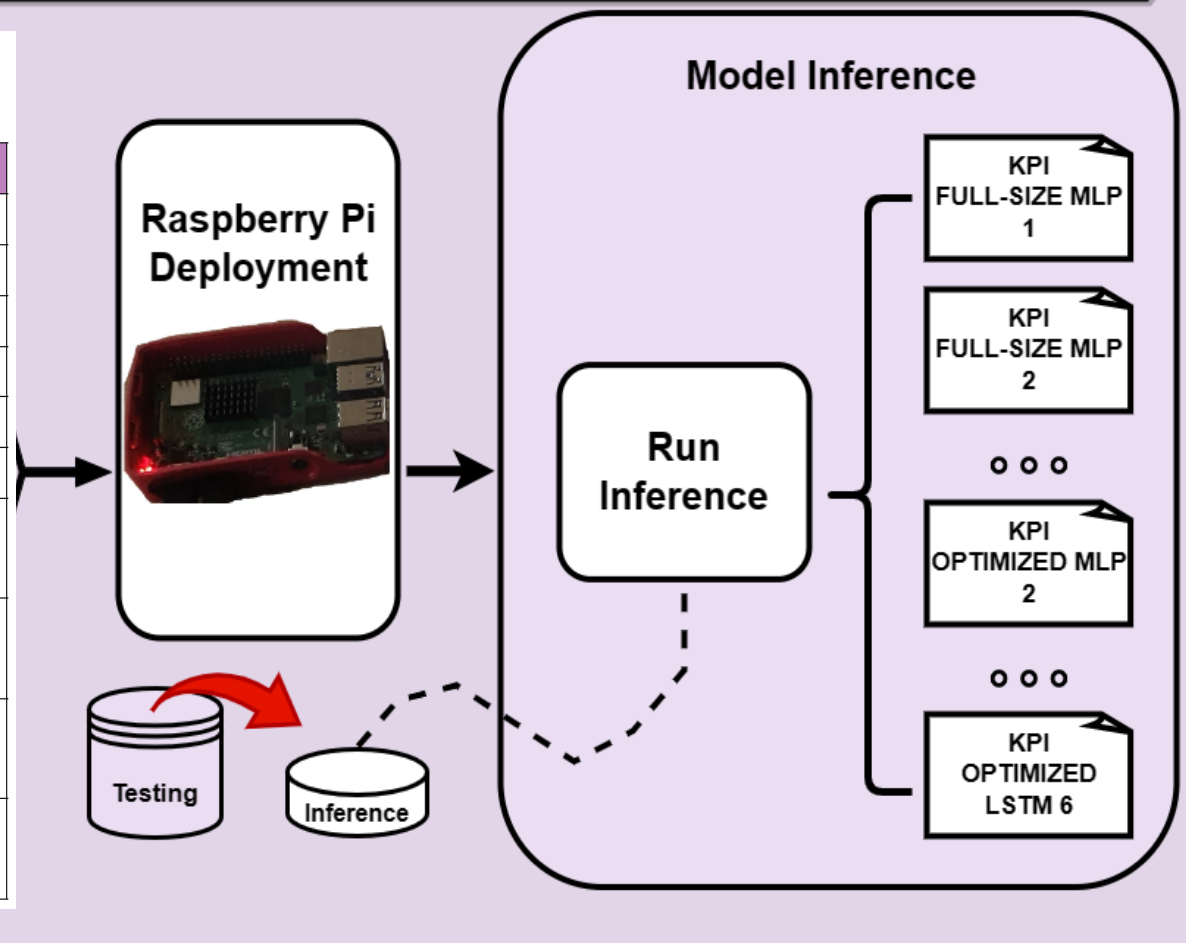
Eksperiment oppsett (II)

Model Optimization

Model Deployment and Inference

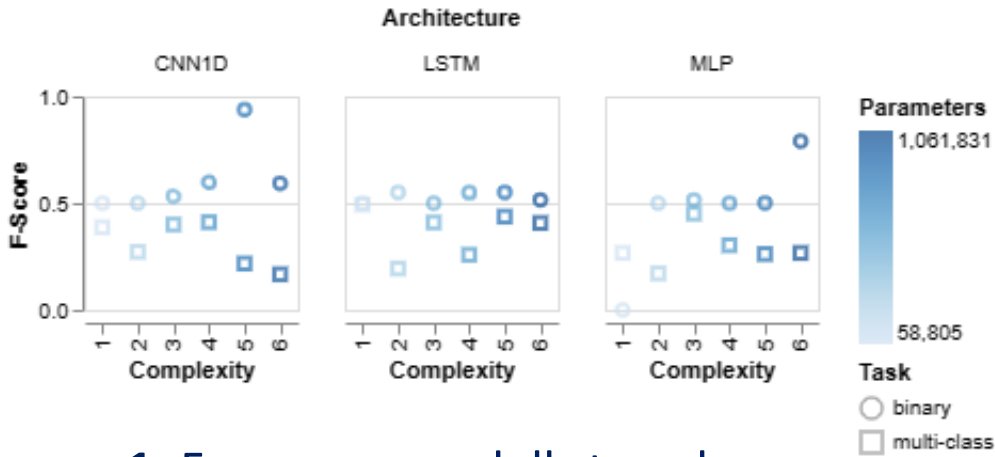
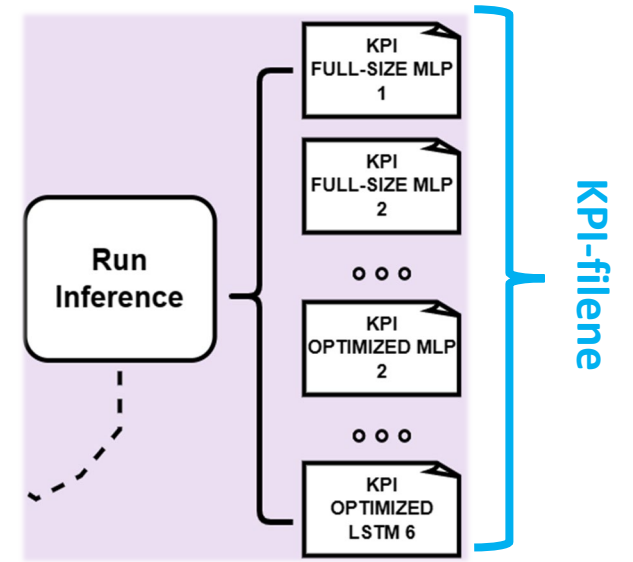
Table 6: Model optimization descriptions.

Optimization ID	Optimization Short Description
baseline	Dynamic range quantization
int8	Quantization integer only (8-bit)
int16	Quantization integer with float fallback (16-bit)
int32	Quantization integer with float fallback (32-bit)
float16	Quantization float (16-bit)
prun	Magnitude-based weight pruning
pruint8	Magnitude-based weight pruning + Quantization integer only (8-bit)
pruint16	Magnitude-based weight pruning + Quantization integer with float fallback (16-bit)
pruint32	Magnitude-based weight pruning + Quantization integer with float fallback (32-bit)
prunfloat16	Magnitude-based weight pruning + Quantization float (16-bit)

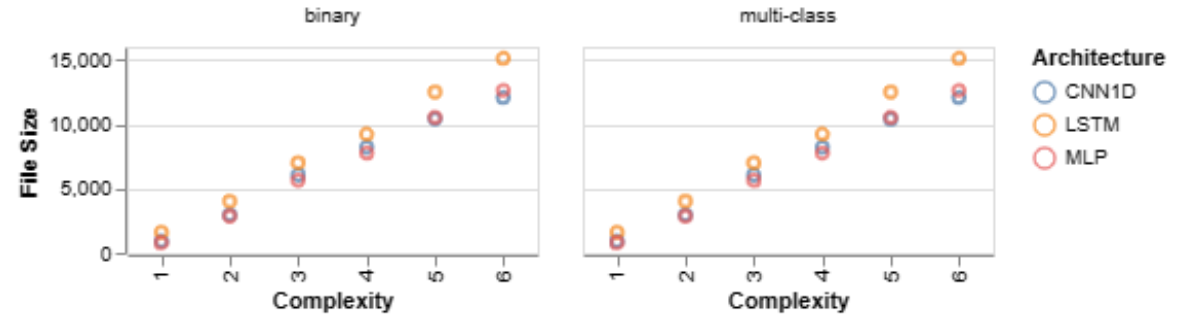


Analyse

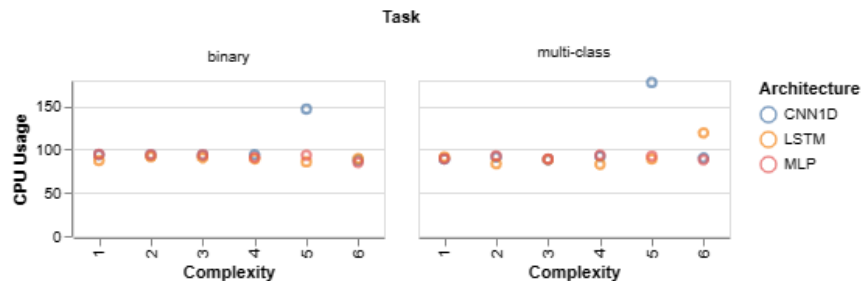
- Vi bruker **KPI-filene** til å analysere:
- (1) Ressurskravene til de ikke-optimiserte modellene.



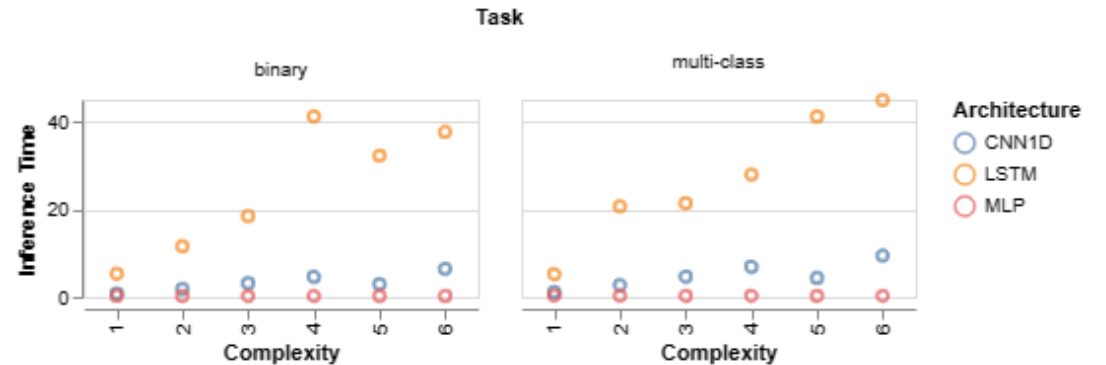
1. F-score vs modell størrelse.



2. Filstørrelse vs modellstørrelse.



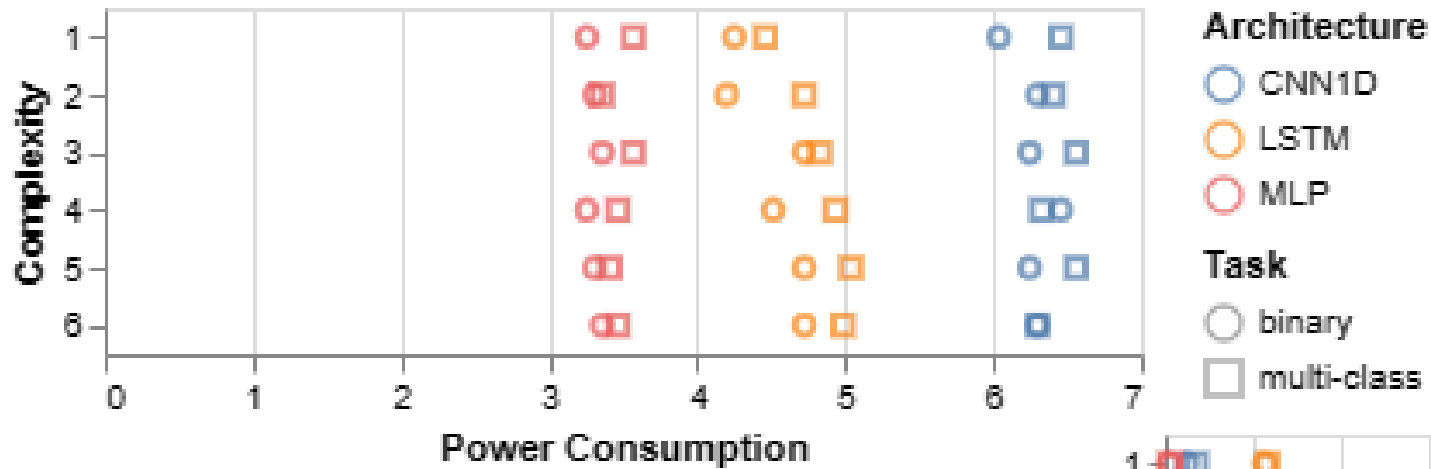
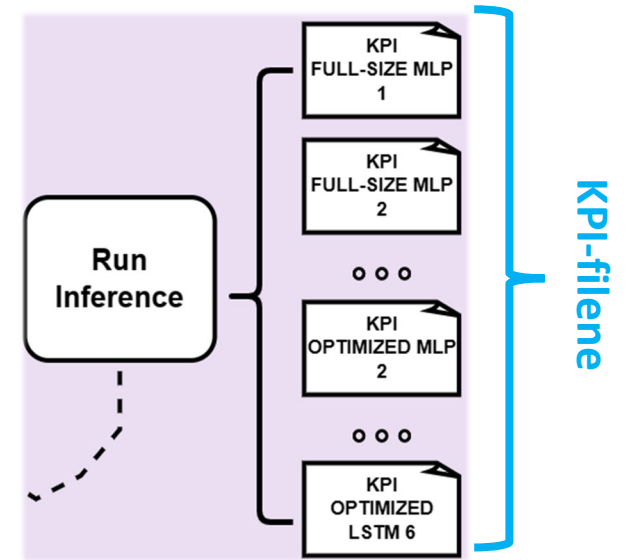
3. CPU-bruk vs modellstørrelse.



4. Inferenstid vs modellstørrelse.

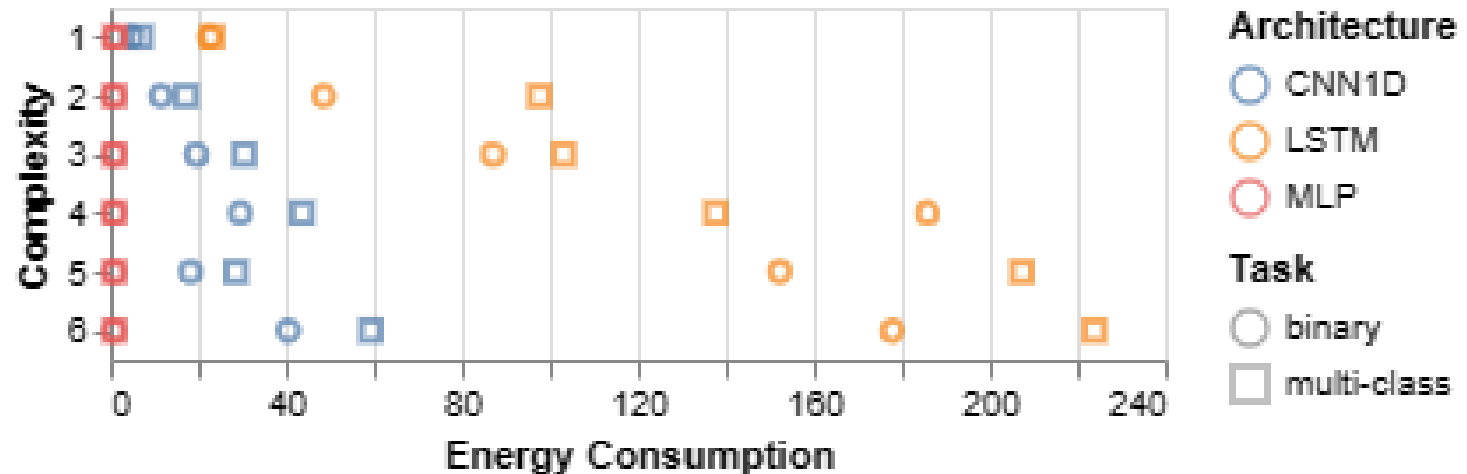
Analyse

- Vi bruker **KPI-filene** til å analysere:
- (1) Ressurskravene til de ikke-optimiserte modellene.



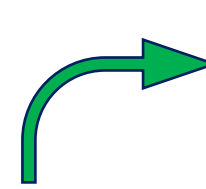
5. Strømforbruk vs modellstørrelse.

6. Energiforbruk vs modellstørrelse.



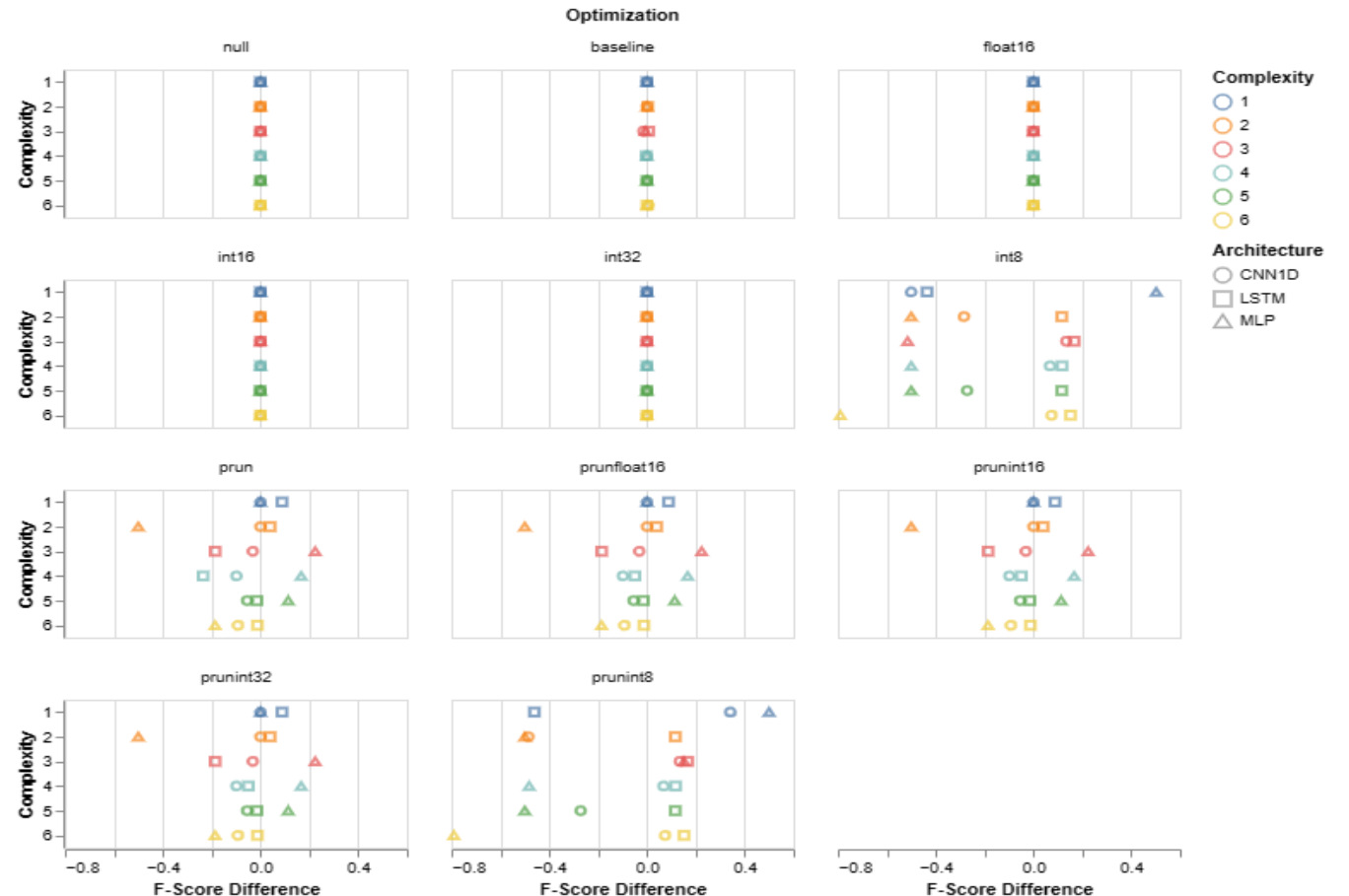
Analyse

- Vi bruker **KPI-filene** til å analysere:
- (2) Om de optimiserte modellene oppfyller **suksesskriteriene**.



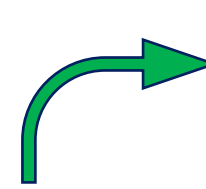
- ✓ Redusere ressursutnyttelsen.
- ✓ Bevare nøyaktigheten.

For å sjekke om nøyaktigheten er bevart beregner man differansen i F-score mellom optimiserte og ikke-optimiserte modeller.



Analyse

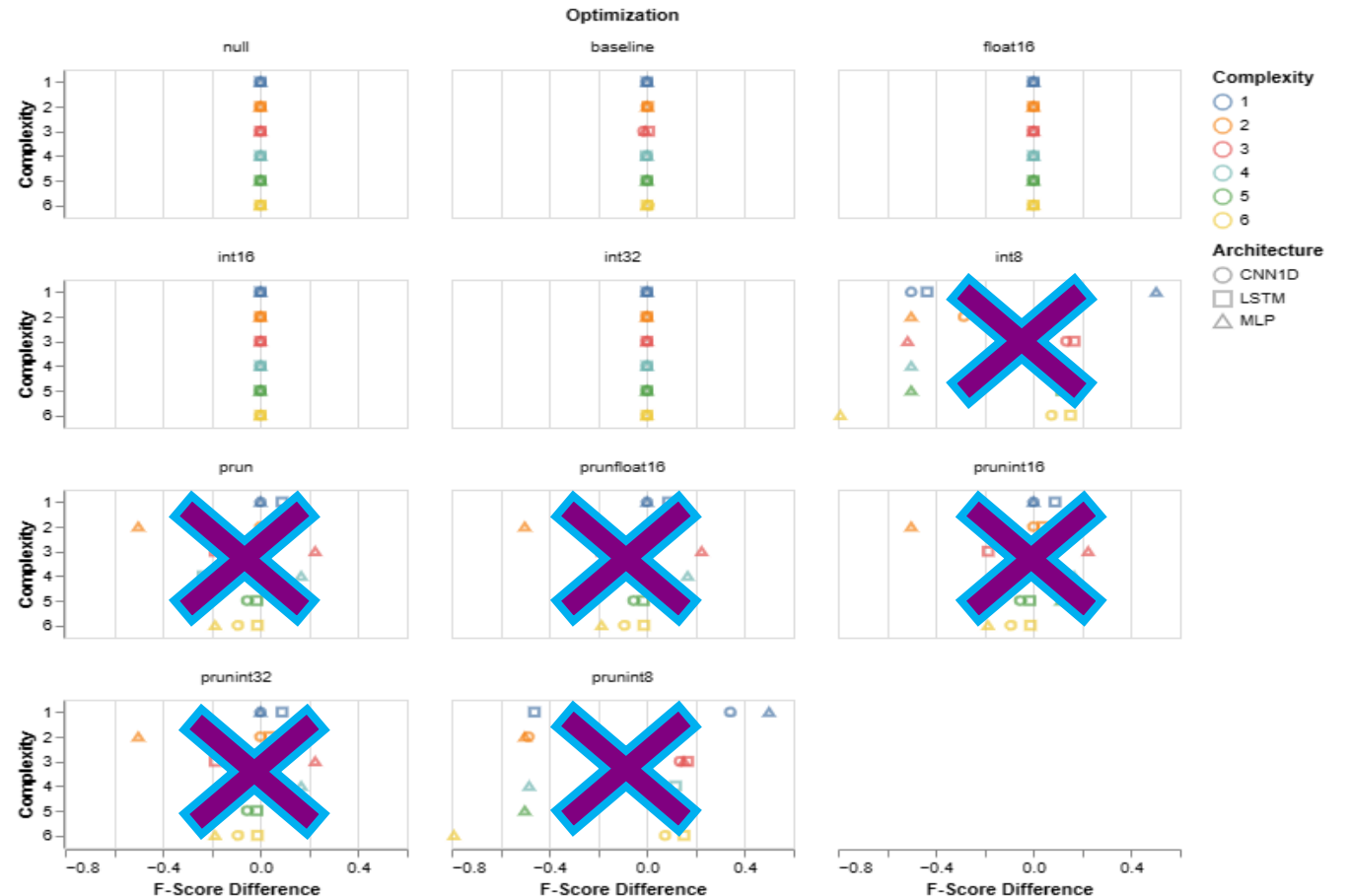
- Vi bruker **KPI-filene** til å analysere:
- (2) Om de optimiserte modellene oppfyller **suksesskriteriene**.



- ✓ Redusere ressursutnyttelsen.
- ✓ Bevare nøyaktigheten.

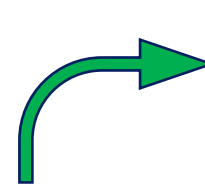
For å sjekke om nøyaktigheten er bevart beregner man differansen i F-score mellom optimiserte og ikke-optimiserte modeller.

Kvantiseringsoptimaliseringer (unntatt int8) oppfyller dette kriteriet!



Analyse

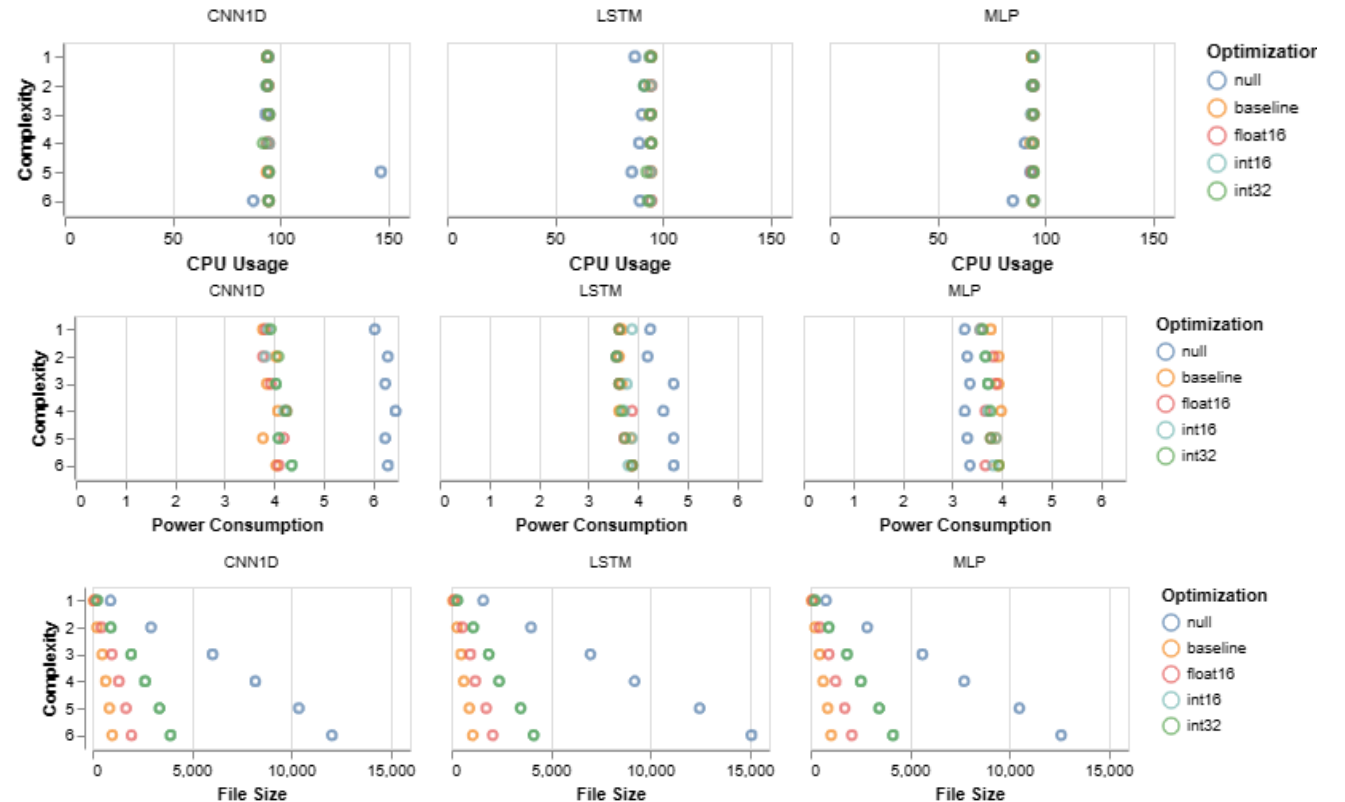
- Vi bruker **KPI-filene** til å analysere:
- (2) Om de optimiserte modellene oppfyller **suksesskriteriene**.



- ✓ Redusere ressursutnyttelsen.
- ✓ Bevare nøyaktigheten.

For å sjekke om ressursbruken er redusert, sjekker vi alle ressurser:

- CPU-bruk vs modellstørrelse.
- Strømforbruk vs modellstørrelse.
- Minneforbruk vs modellstørrelse.

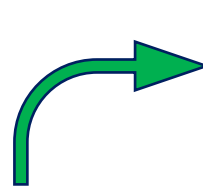




SINTEF

Analyse

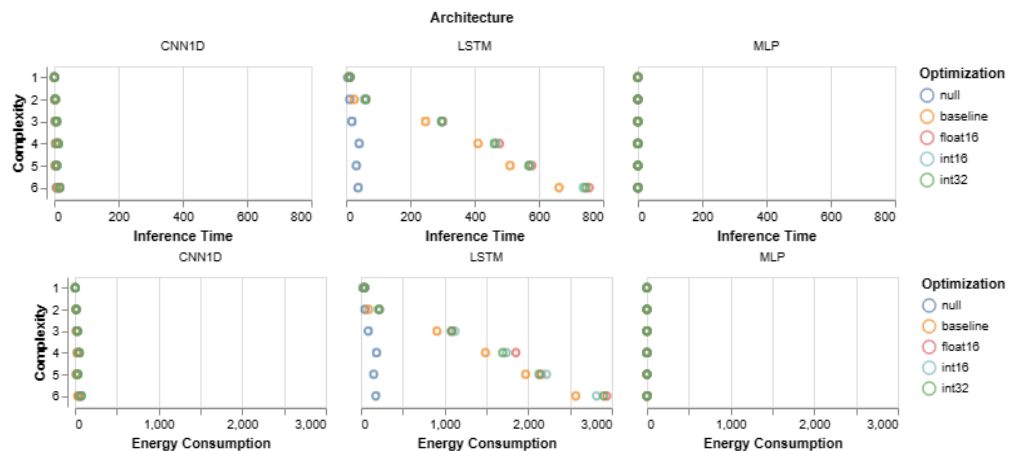
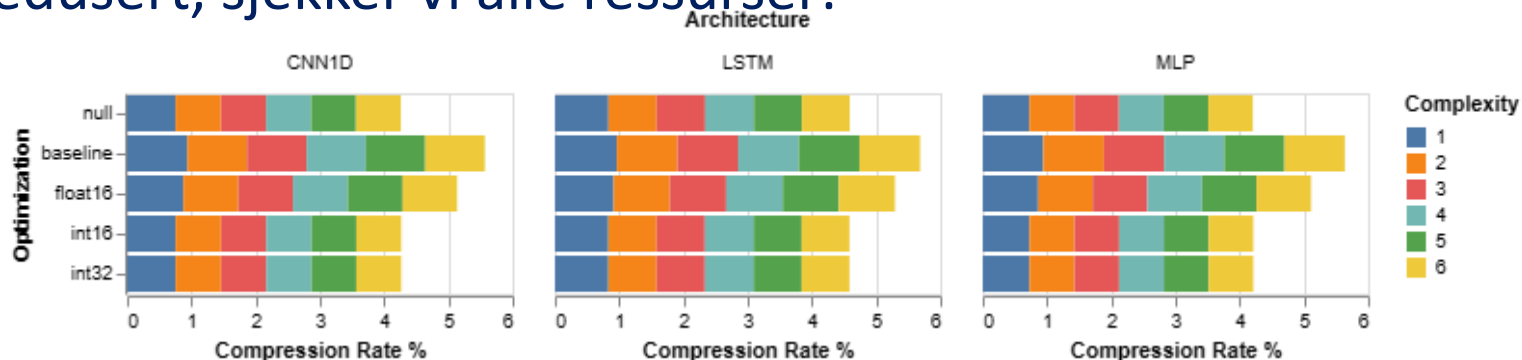
- Vi bruker **KPI-filene** til å analysere:
- (2) Om de optimiserte modellene oppfyller **suksesskriteriene**.



- ✓ Redusere ressursutnyttelsen.
- ✓ Bevare nøyaktigheten.

For å sjekke om ressursbruken er redusert, sjekker vi alle ressurser:

- Komprimeringsgrad vs modellstørrelse.
- Inferenstid/energiforbruk vs modellstørrelse.





SINTEF

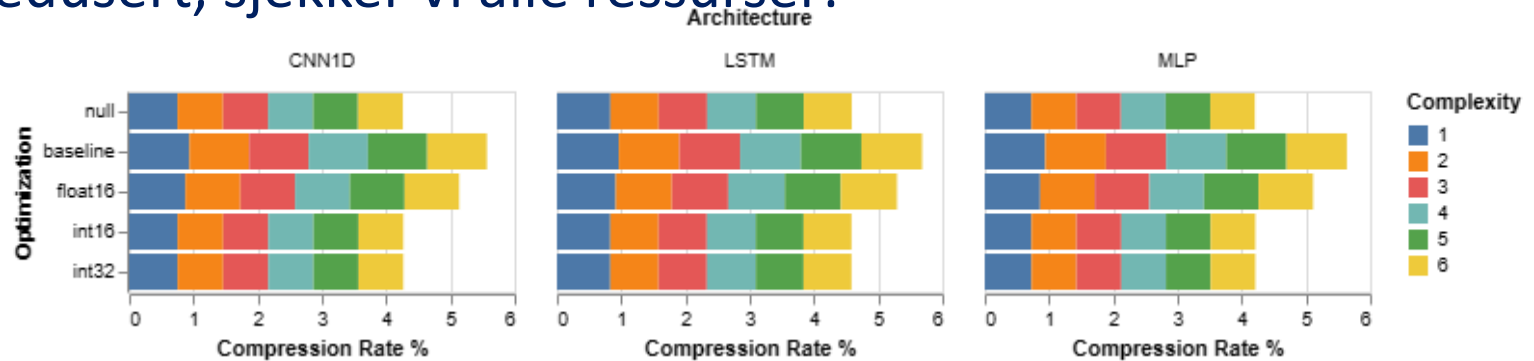
Analyse

- Vi bruker **KPI-filene** til å analysere:
- (2) Om de optimiserte modellene oppfyller **suksesskriteriene**.

✓ Redusere ressursutnyttelsen.
 ✓ Bevare nøyaktigheten.

For å sjekke om ressursbruken er redusert, sjekker vi alle ressurser:

- Komprimeringsgrad vs modellstørrelse.



- Inferenstid/energiforbruk vs modellstørrelse.



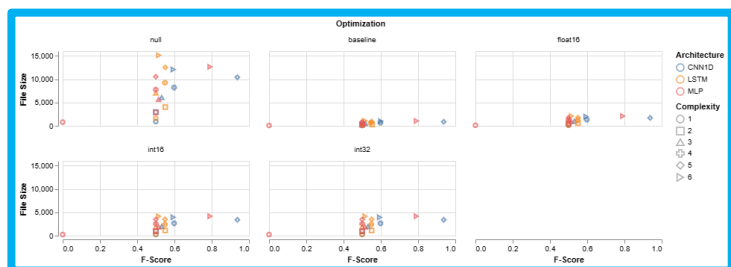
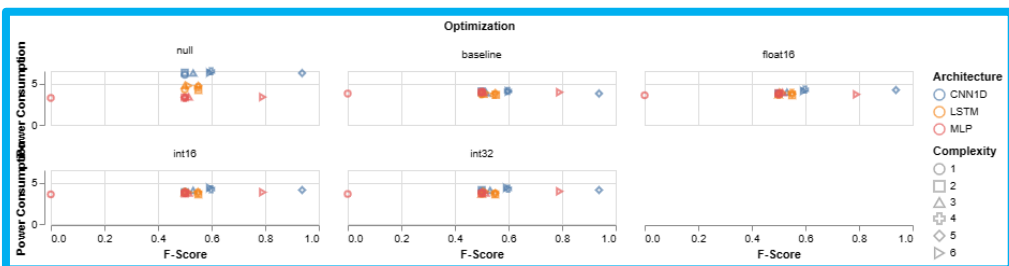
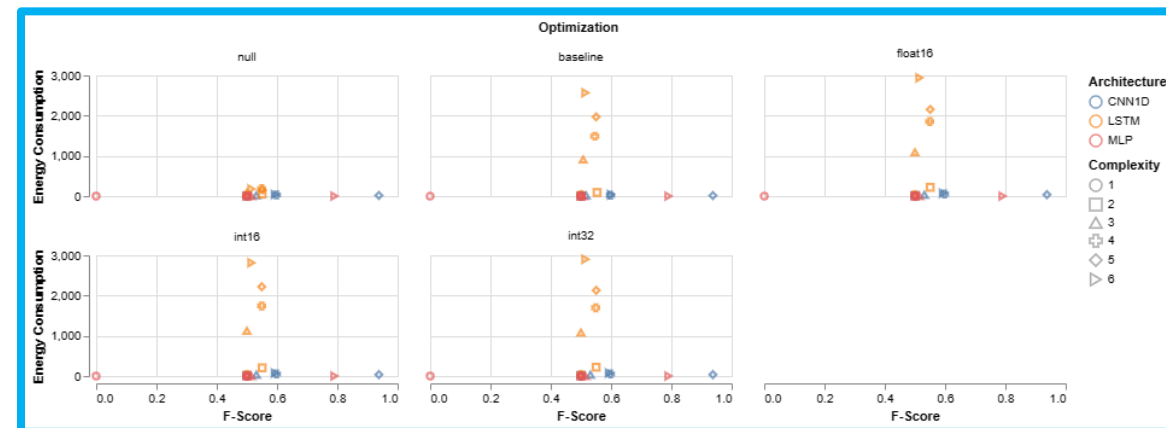
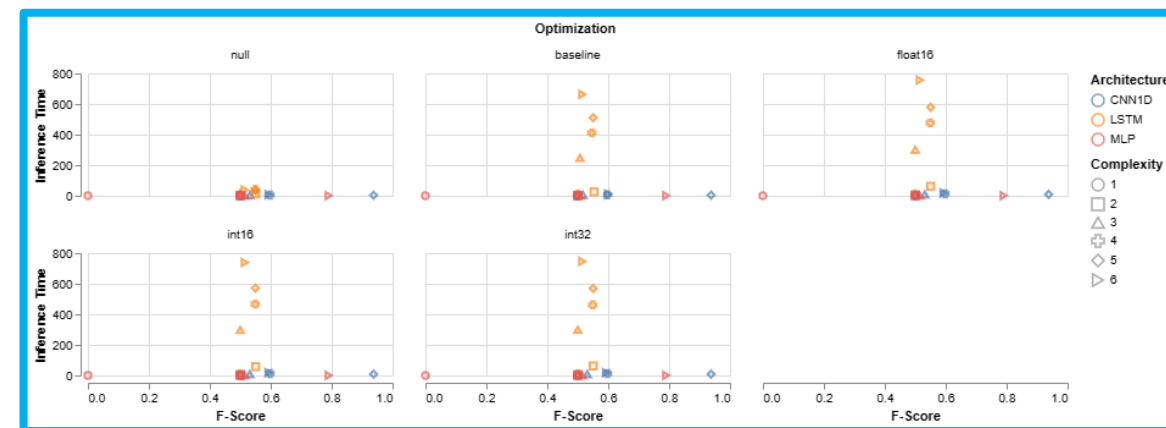
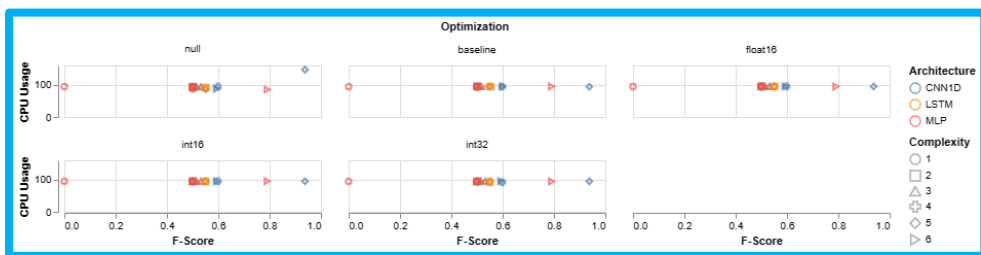
OBS! Ikke-optimiserte LSTM-modeller presterer best!



SINTEF

Analyse

- Vi bruker **KPI-filene** til å analysere:
- (3) Ressurs vs nøyaktighet **avveining**: hvilke optimaliseringer oppnår en bedre **balanse**?





Resultater

(1) → Ulike arkitekturer har ulike behov for optimalisering.



Resultater

- (1) → Ulike arkitekturer har ulike behov for optimalisering.
- (2) → Vi finner at ingen av optimaliseringene fra kandidatsettet oppfyller kriteriene.



Resultater

(1) → Ulike arkitekturer har ulike behov for optimalisering.

(2) → Vi finner at ingen av optimaliseringene fra kandidatsettet oppfyller kriteriene.



(3) → **MEN!**

- baseline og float16 har den beste avveiningen når det gjelder minneforbruk.
- baseline har den beste avveiningen når det gjelder strømforbruk.
- alle optimiseringer viser svært liknende resultater for resten av ressursene.



Resultater

(1) → Ulike arkitekturer har ulike behov for optimalisering.

(2) → Vi finner at ingen av optimaliseringene fra kandidatsettet oppfyller kriteriene.



(3) → **MEN!**

- baseline og float16 har den beste avveiningen når det gjelder minneforbruk.
- baseline har den beste avveiningen når det gjelder strømforbruk.
- alle optimiseringer viser svært liknende resultater for resten av ressursene.



→ Rammeverket lykkes med å eliminere optimiseringsteknikk som ikke oppfyller suksesskriteriene og fange avveininger mellom nøyaktighet og ressurser.

Resultater

(1) → Ulike arkitekturer har ulike behov for optimalisering.

(2) → Vi finner at ingen av optimaliseringene fra kandidatsettet oppfyller kriteriene.



(3) → **MEN!**

- baseline og float16 har den beste avveiningen når det gjelder minneforbruk.
- baseline har den beste avveiningen når det gjelder strømforbruk.
- alle optimiseringer viser svært liknende resultater for resten av ressursene.



→ Rammeverket lykkes med å eliminere optimiseringsteknikk som ikke oppfyller suksesskriteriene og fange avveininger mellom nøyaktighet og ressurser.

→ Eksperimentresultater viser at MLP og CNN modeller kan optimaliseres ved hjelp av **baseline, int16, int32 og float16**, og samtidig oppfylle begge to kriteriene.



Resultater

(1) → Ulike arkitekturer har ulike behov for optimalisering.

(2) → Vi finner at ingen av optimaliseringene fra kandidatsettet oppfyller kriteriene.



(3) → **MEN!**

- baseline og float16 har den beste avveiningen når det gjelder minneforbruk.
- baseline har den beste avveiningen når det gjelder strømforbruk.
- alle optimiseringer viser svært liknende resultater for resten av ressursene.



→ Rammeverket lykkes med å eliminere optimiseringsteknikk som ikke oppfyller suksesskriteriene og fange avveininger mellom nøyaktighet og ressurser.

→ Eksperimentresultater viser at MLP og CNN modeller kan optimaliseres ved hjelp av **baseline, int16, int32** og **float16**, og samtidig oppfylle begge to kriteriene.

→ Pruning kan i noen tilfeller forbedre nøyaktighetsgraden til en modell.



Resultater

(1) → Ulike arkitekturer har ulike behov for optimalisering.

(2) → Vi finner at ingen av optimaliseringene fra kandidatsettet oppfyller kriteriene.



(3) → **MEN!**

- baseline og float16 har den beste avveiningen når det gjelder minneforbruk.
- baseline har den beste avveiningen når det gjelder strømforbruk.
- alle optimiseringer viser svært liknende resultater for resten av ressursene.



→ Rammeverket lykkes med å eliminere optimiseringsteknikk som ikke oppfyller suksesskriteriene og fange avveininger mellom nøyaktighet og ressurser.

→ Eksperimentresultater viser at MLP og CNN modeller kan optimaliseres ved hjelp av **baseline, int16, int32** og **float16**, og samtidig oppfylle begge to kriteriene.

→ Pruning kan i noen tilfeller forbedre nøyaktighetsgraden til en modell.

→ Rammeverket kan forbedres!



SINTEF

Referanser

Masteroppgave:

Videsjorden, A. N. (2023). *Optimizing deep learning inference for resource-constricted platforms*

- [1] Zheng, T., Ardolino, M., Bacchetti, A., & Perona, M. (2021). The applications of Industry 4.0 technologies in manufacturing context: a systematic literature review. *International Journal of Production Research*, 59(6), 1922-1954.
- [2] Rai, R., Tiwari, M. K., Ivanov, D., & Dolgui, A. (2021). Machine learning in manufacturing and industry 4.0 applications. *International Journal of Production Research*, 59(16), 4773-4778.
- [3] Siccoli, A., de Wispelaere, M. P., Schröder, M. L., & Staartjes, V. E. (2019). Machine learning–based preoperative predictive analytics for lumbar spinal stenosis. *Neurosurgical Focus*, 46(5), E5.
- [4] Bannigan, P., Aldeghi, M., Bao, Z., Häse, F., Aspuru-Guzik, A., & Allen, C. (2021). Machine learning directed drug formulation development. *Advanced Drug Delivery Reviews*, 175, 113806.
- [5] Bai, Y., Guan, Y., & Ng, W. F. (2020, September). Fatigue assessment using ECG and actigraphy sensors. In *Proceedings of the 2020 ACM International Symposium on Wearable Computers* (pp. 12-16).
- [6] Garcia-Ceja, E., Osmani, V., & Mayora, O. (2015). Automatic stress detection in working environments from smartphones' accelerometer data: a first step. *IEEE journal of biomedical and health informatics*, 20(4), 1053-1060.
- [7] Little, B., Alshabrawy, O., Stow, D., Ferrier, I. N., McNaney, R., Jackson, D. G., ... & O'Brien, J. T. (2021). Deep learning-based automated speech detection as a marker of social functioning in late-life depression. *Psychological Medicine*, 51(9), 1441-1450.
- [8] Tang, C. S., & Veelenturf, L. P. (2019). The strategic role of logistics in the industry 4.0 era. *Transportation Research Part E: Logistics and Transportation Review*, 129, 1-11.
- [9] Bouman, P., Agatz, N., & Schmidt, M. (2018). Dynamic programming approaches for the traveling salesman problem with drone. *Networks*, 72(4), 528-542.
- [10] Sajjad, M., Irfan, M., Muhammad, K., Del Ser, J., Sanchez-Medina, J., Andreev, S., ... & Lee, J. W. (2020). An efficient and scalable simulation model for autonomous vehicles with economical hardware. *IEEE transactions on intelligent transportation systems*, 22(3), 1718-1732.
- [11] Nielsen, M. A. (2015). *Neural networks and deep learning* (Vol. 25, pp. 15-24). San Francisco, CA, USA: Determination press.
- [12] Coelho, Y. L., dos Santos, F. D. A. S., Frizzera-Neto, A., & Bastos-Filho, T. F. (2021). A lightweight framework for human activity recognition on wearable devices. *IEEE Sensors Journal*, 21(21), 24471-24481.
- [13] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT press
- [14] Buyya, R., Broberg, J., & Goscinski, A. M. (Eds.). (2010). *Cloud computing: Principles and paradigms*. John Wiley & Sons.
- [15] Satyanarayanan, M., Bahl, P., Caceres, R., & Davies, N. (2009). The case for vm-based cloudlets in mobile computing. *IEEE pervasive Computing*, 8(4), 14-23.
- [16] Parikh, S., Dave, D., Patel, R., & Doshi, N. (2019). Security and privacy issues in cloud, fog and edge computing. *Procedia Computer Science*, 160, 734-739.
- [17] Liu, D., Kong, H., Luo, X., Liu, W., & Subramaniam, R. (2022). Bringing AI to edge: From deep learning's perspective. *Neurocomputing*, 485, 297-320.



SINTEF

Takk!
Q & A