Transformer Networks Intuition, Theory, and Real-world Deployment

Daniel Romero

Dept. of Information and Communication Technology
University of Agder, Norway



Acknowledgements:

Pham Q. Viet

Research Council of Norway (IKTPLUSS grant 311994)



The Al wave...

Social focus on AI mainly due to advances in large language models (LLMs)

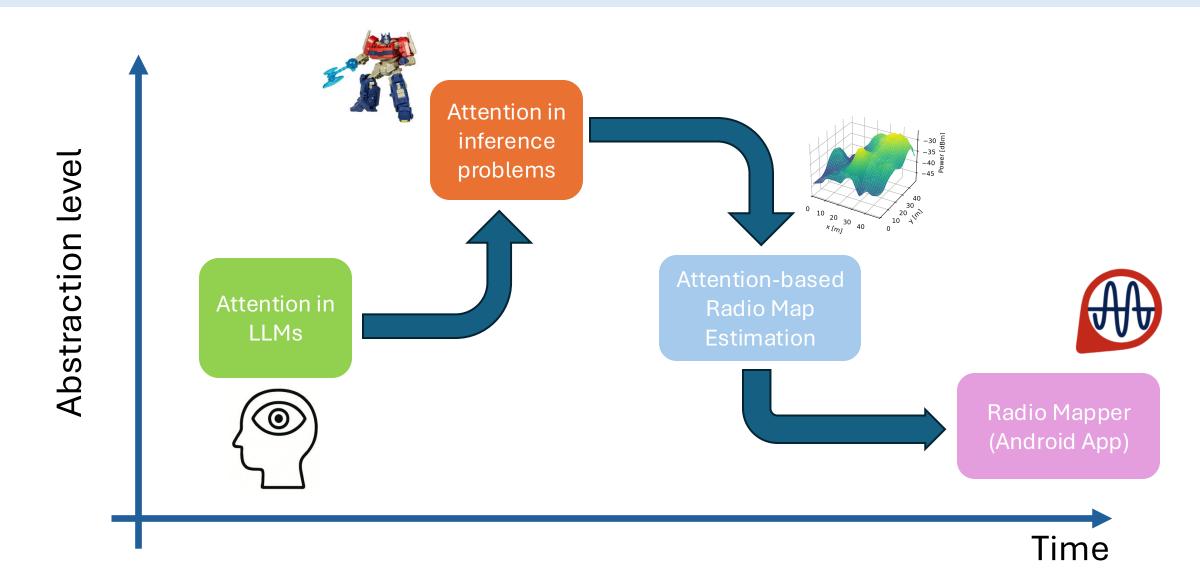


- Goal of this talk:
 - > Understand the intuition and theoretical principles in

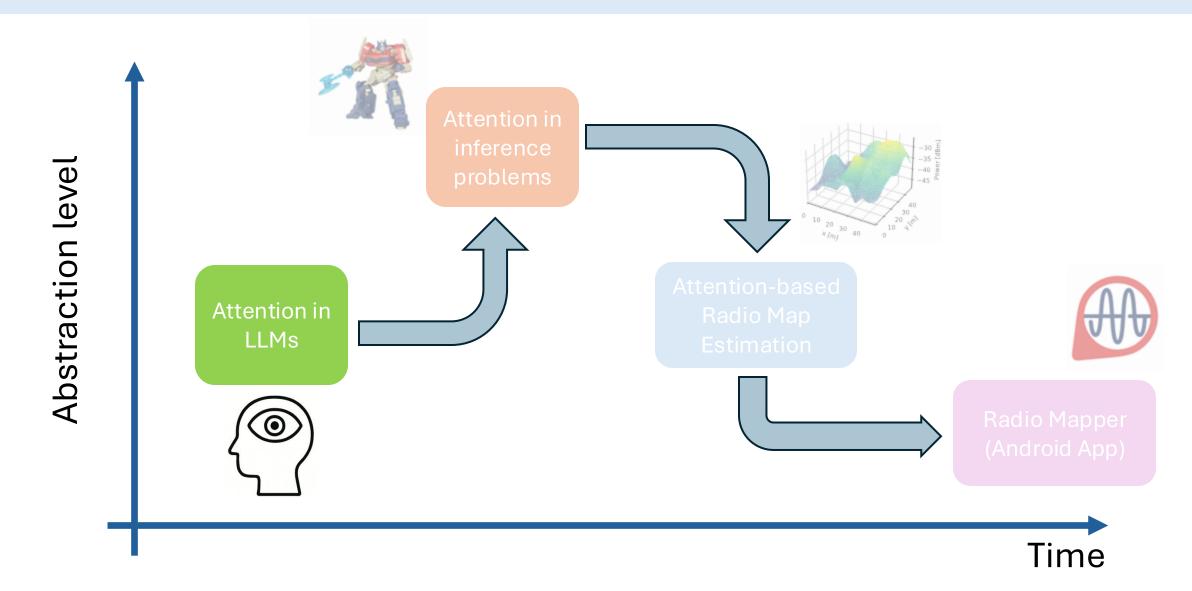
Attention

> Learn to apply these ideas to inference problems beyond LLMs

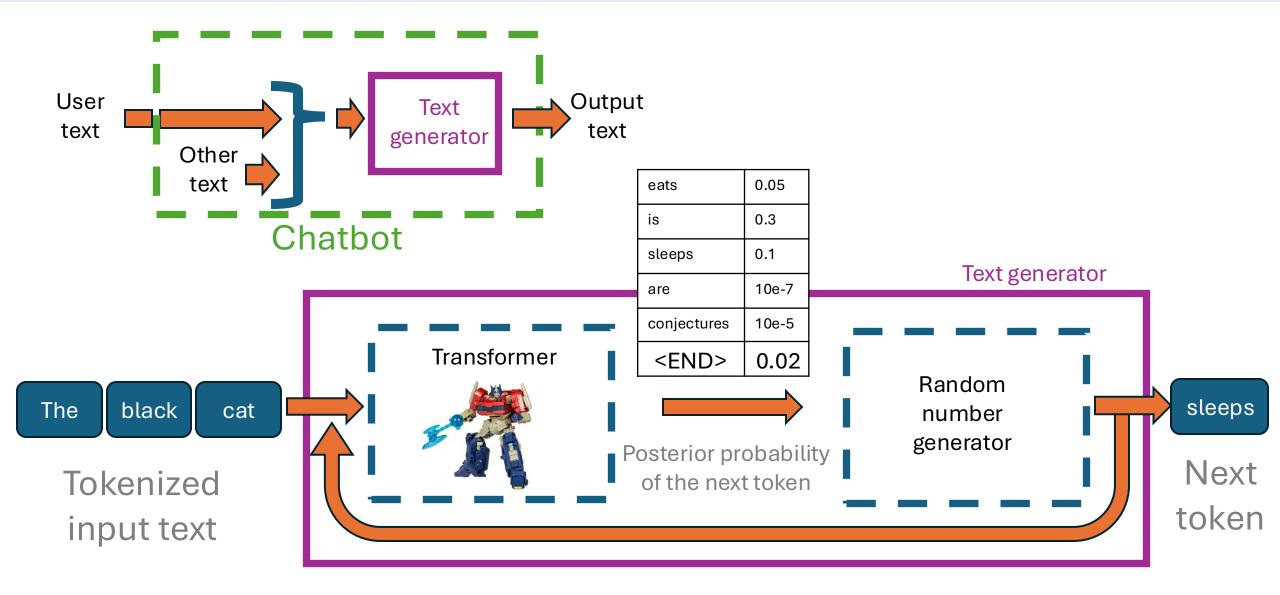
Outline of this talk



Outline of this talk



Autoregressive Text Generation



Append the generated token to the input and repeat until <END> is generated

Attention in NLP

- Combining words results in richer concepts
 - Word relations based on syntax
 - The house is white.
 - The car near the house is white.
 - Word relations based on semantics
 - The cat ate the food because it was hungry.
 - The cat ate the food because it was warm.





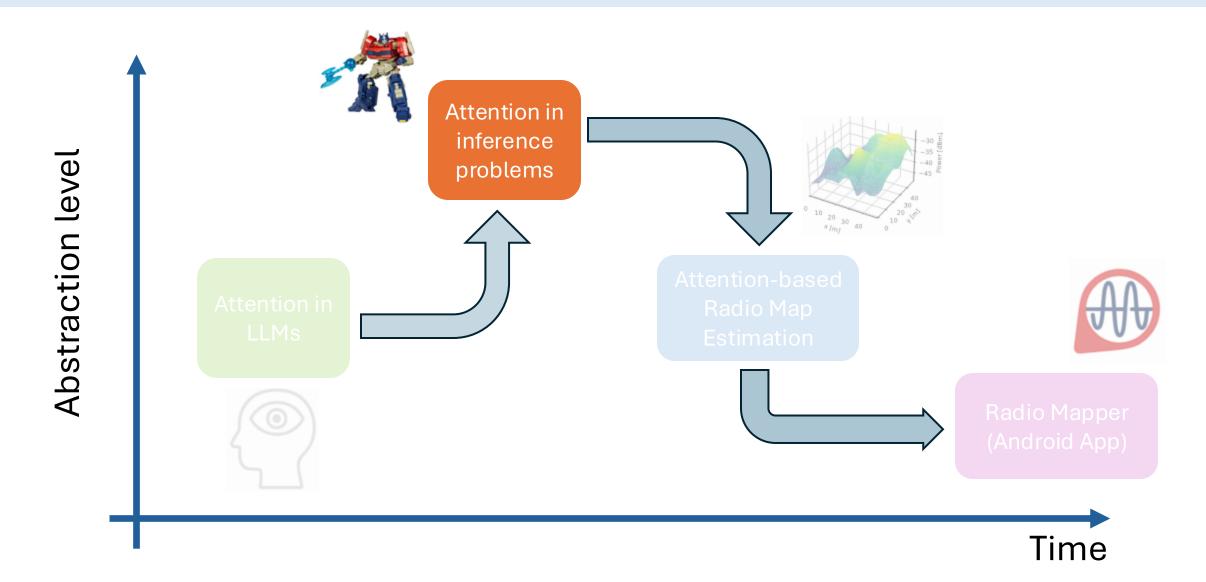


Each word

"pays more attention" to some words than to others.

"enriches the meaning" of related words.

Outline of this talk



Attention in Inference Problems

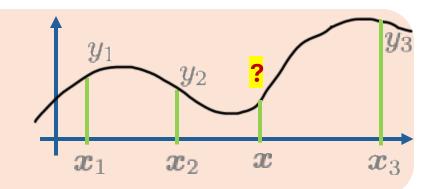
Autoregressive text generation

Given
$$[u[1],\ldots,u[T]]$$
, obtain the posterior prob. $p\in[0,1]^{N_{\mathrm{U}}}$ of the next token $\mathcal{U}:=\{u^{(1)},\ldots,u^{(N_{\mathrm{U}})}\}$

Time series forecasting

Given
$$m{s}[au-T],\ldots,m{s}[au-1]$$
 , predict $m{s}[au]$ $\in \mathbb{R}^{D_{m{x}}}$



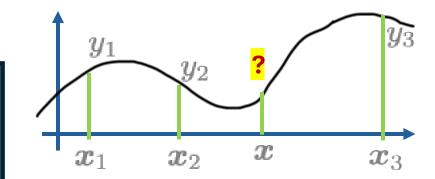


Note that the number of inputs (N, T) is arbitrary!

Weighted Averages for Function Estimation

Given
$$\{({m x}_n,y_n)\}_{n=1}^N$$
 and ${m x}$, estimate $y:=f({m x})$

• Estimator (general form):
$$\hat{y} = \frac{\sum_{n=1}^{N} \alpha_n y_n}{\sum_{n=1}^{N} \alpha_n}$$



Prominent special cases

> K-NN	$\alpha_n = \begin{cases} 1 & \text{if } \boldsymbol{x}_n \text{ a NN of } \boldsymbol{x}, \\ 0 & \text{otherwise.} \end{cases}$	
> Radius-NN	$\alpha_n = \begin{cases} 1 & \text{if } \ \boldsymbol{x} - \boldsymbol{x}_n\ _2 \le R, \\ 0 & \text{otherwise.} \end{cases}$:= $\alpha(\boldsymbol{x}_n, \boldsymbol{x})$	
Shepard's method	$\alpha_n = \frac{1}{\ \boldsymbol{x} - \boldsymbol{x}_n\ _2^{\gamma}} := \alpha(\boldsymbol{x}_n, \boldsymbol{x})$	
Nadaraya-Watson estimator	$lpha_n = lpha(oldsymbol{x}_n, oldsymbol{x})$ a kernel, e.g. Gaussian	

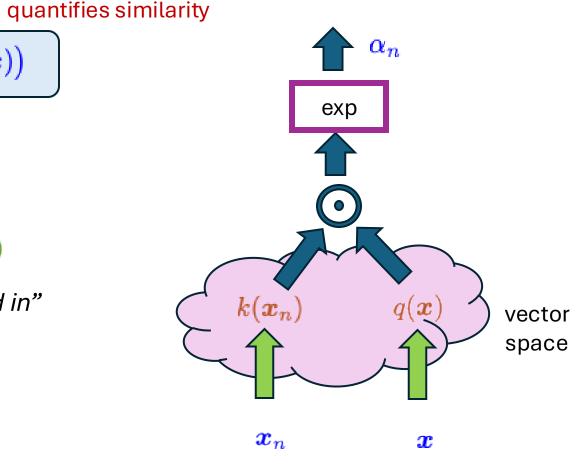
 $\alpha_n \rightarrow$ Quantify how relevant the n-th observation is \rightarrow Attention weights

Inner-Product Attention

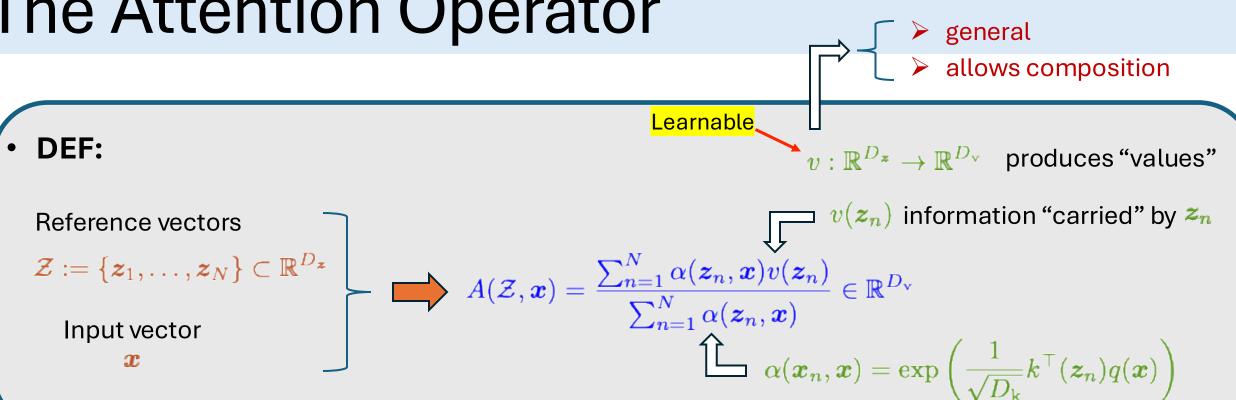
• Given by
$$\alpha_n = \alpha({m x}_n, {m x}) = \exp\left(k^{ op}({m x}_n) \cdot q({m x})\right)$$

- As a soft-max: $\hat{y} = \frac{\sum_{n=1}^{N} \alpha_n y_n}{\sum_{n=1}^{N} \alpha_n} = \sum_{n=1}^{N} \bar{\alpha}_n y_n$

$$ar{lpha}(oldsymbol{x}_n,oldsymbol{x}) = rac{\exp\left(k^{ op}(oldsymbol{x}_n)q(oldsymbol{x})
ight)}{\sum_{n'=1}^N \exp\left(k^{ op}(oldsymbol{x}_{n'})q(oldsymbol{x})
ight)} \quad pprox \delta[n-n_{ ext{max}}]$$
 if the $n_{ ext{max}}$ -th term dominates



The Attention Operator

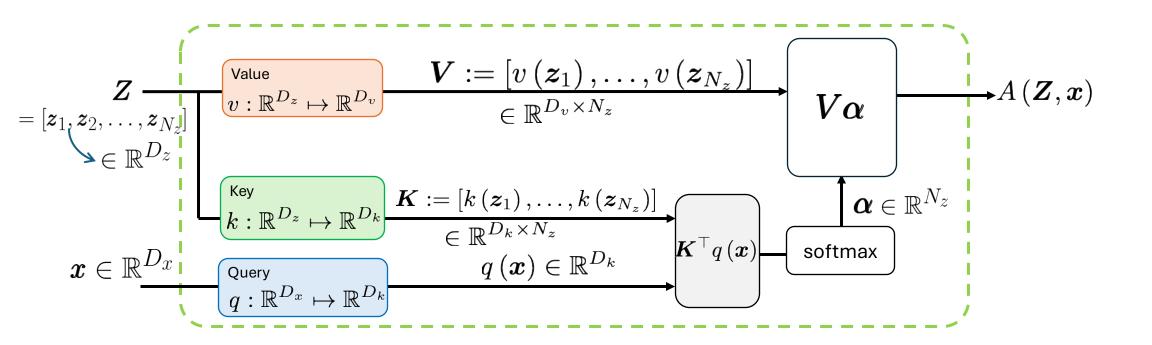


Intuition: This operator enriches the information in **w** with the information in the reference vectors

These functions are typically **linear**: $k(\boldsymbol{z}_n) = \boldsymbol{W}_k \boldsymbol{z}_n, \quad q(\boldsymbol{x}) = \boldsymbol{W}_q \boldsymbol{x}, \quad v(\boldsymbol{z}_n) = \boldsymbol{W}_v \boldsymbol{z}_n$ Trainable parameters: $\Theta := (W_k, W_q, W_v)$

The Attention Operator

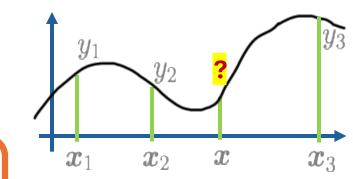
$$A\left(\boldsymbol{Z}, \boldsymbol{x}\right) = \boldsymbol{V} \operatorname{softmax}\left(\boldsymbol{K}^{\top} q\left(\boldsymbol{x}\right)\right)$$



Attention Operator in Function Estimation

• The attention operator generalizes the function estimator seen earlier:

$$A(\mathbf{Z}, \mathbf{x}) = \frac{\sum_{n=1}^{N} \alpha(\mathbf{z}_n, \mathbf{x}) v(\mathbf{z}_n)}{\sum_{n=1}^{N} \alpha(\mathbf{z}_n, \mathbf{x})} \qquad \qquad \hat{y} = \frac{\sum_{n=1}^{N} \alpha_n y_n}{\sum_{n=1}^{N} \alpha_n}$$



$$m{z}_n = \left[egin{array}{c} y_n \ m{x}_n \end{array}
ight] \quad k(m{z}_n) = [m{0},ar{m{W}}_k]m{z}_n \quad v(m{z}_n) = [1,m{0}^ op]m{z}_n = y_n$$

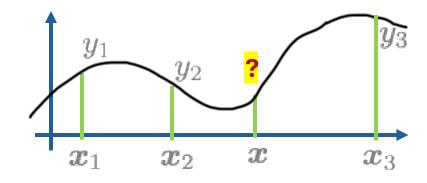
In practice, if the problem is invariant to translations

$$extstyle m{z}_n = egin{bmatrix} y_n \ m{x}_n - m{x} \end{bmatrix}$$
 $\hat{y} = A(\mathcal{Z}, m{z}_N)$ (the estimator above is not invariant)

Training (Toy) Example

• Consider a dataset for the function estimation problem

$$\{(\boldsymbol{x}_n^{(j)},y_n^{(j)}):\quad n=0,\ldots,N, j=1,\ldots,J\}$$
 number of realizations



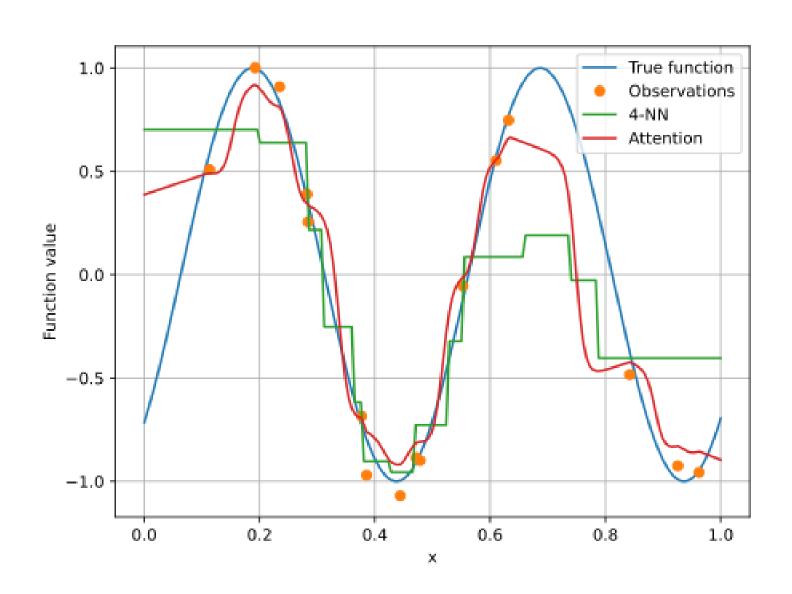
Simple MSE loss (1 target/realization)

$$\mathcal{L}(\Theta) = rac{1}{J} \sum_{j=1}^J \left(oldsymbol{y}_0^{(j)} - A(oldsymbol{\mathcal{Z}}^{(j)}, oldsymbol{z}_0^{(j)})
ight)^2 \qquad \qquad oldsymbol{\mathcal{Z}}^{(j)} := \left\{ oldsymbol{z}_1^{(j)}, \dots, oldsymbol{z}_N^{(j)}
ight\}$$

- Multiple targets per realization also possible
- Training means finding a local optimum of

$$\underset{\Theta}{\operatorname{minimize}} \ \mathcal{L}(\Theta)$$

Training Example



Attention Heads

• Cross-attention head

$$\mathcal{Z}:=\{m{z}_1,\ldots,m{z}_N\}\subset\mathbb{R}^{D_{m{z}}}$$
 order not important $A(\mathcal{Z},m{X}):=[A(\mathcal{Z},m{x}_1),\ldots,A(\mathcal{Z},m{x}_M)]\in\mathbb{R}^{D_{m{v}} imes M}$ order matters

Self-attention head

$$oldsymbol{X} := [oldsymbol{x}_1, \dots, oldsymbol{x}_M]$$

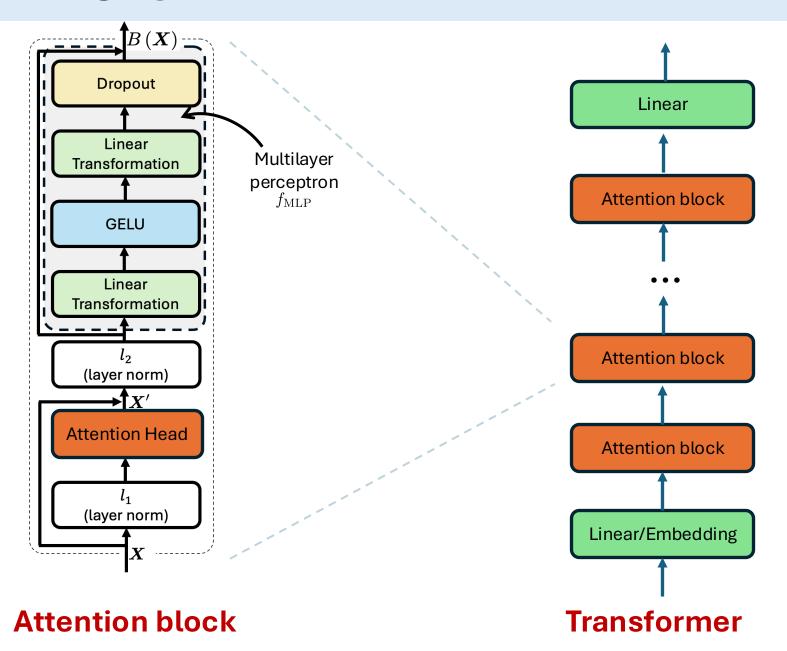
$$A(\boldsymbol{X}) := A([\boldsymbol{x}_1, \dots, \boldsymbol{x}_M]) := A(\{\boldsymbol{x}_1, \dots, \boldsymbol{x}_M\}, [\boldsymbol{x}_1, \dots, \boldsymbol{x}_M])$$

$$= \left[A(\{\boldsymbol{x}_1, \dots, \boldsymbol{x}_M\}, \boldsymbol{x}_1), A(\{\boldsymbol{x}_1, \dots, \boldsymbol{x}_M\}, \boldsymbol{x}_2), \dots, A(\{\boldsymbol{x}_1, \dots, \boldsymbol{x}_M\}, \boldsymbol{x}_M) \right]$$

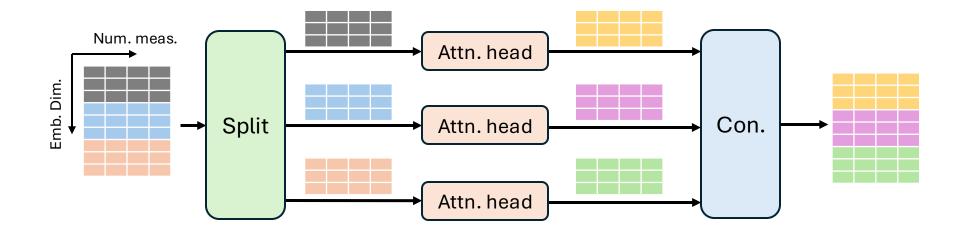
Causal self-attention head

$$oldsymbol{X} := \left[oldsymbol{x}_1, \ldots, oldsymbol{x}_M
ight]$$
 $A_{\mathrm{c}}(oldsymbol{X}) := \left[Aig(\{oldsymbol{x}_1\}, oldsymbol{x}_1ig), Aig(\{oldsymbol{x}_1, oldsymbol{x}_2\}, oldsymbol{x}_2ig), \ldots, Aig(\{oldsymbol{x}_1, \ldots, oldsymbol{x}_M\}, oldsymbol{x}_Mig)
ight]$

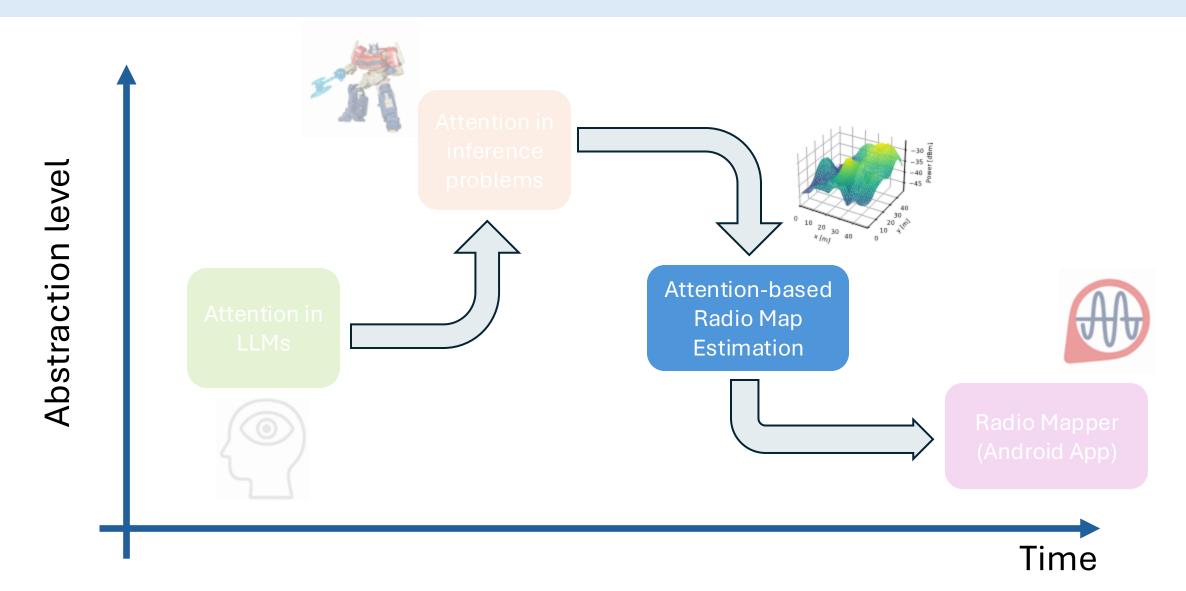
Transformers



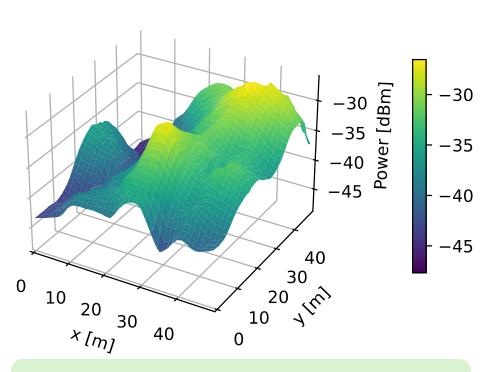
Multi-head Attention



Outline of this talk



Radio Maps



A **radio map** provides a **metric of interest** in a radio communication environment.

Examples:

- Received signal power
- Interference power
- Power spectral density (PSD)
- Electromagnetic absorption
- Channel gain



Overview

Radio map estimation (RME)

Motivation

Determine radio map values where no measurements have been collected

Goal

Interpolative inference based on spatially distributed measurements

Contribution

Transformer-based Estimator (STORM)

Active sensing for RME

Motivation

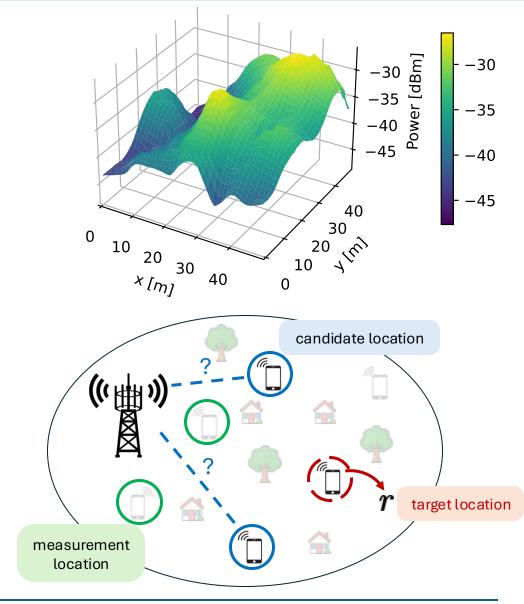
Minimization of drive tests (MDT)

Goal

Choose where to measure next

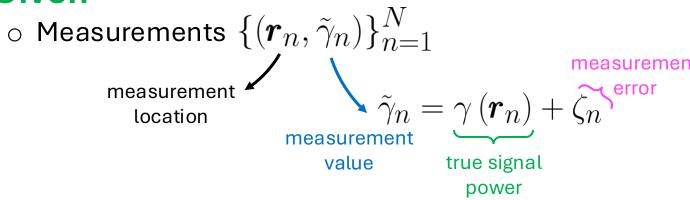
Contribution

Extended transformer-based network



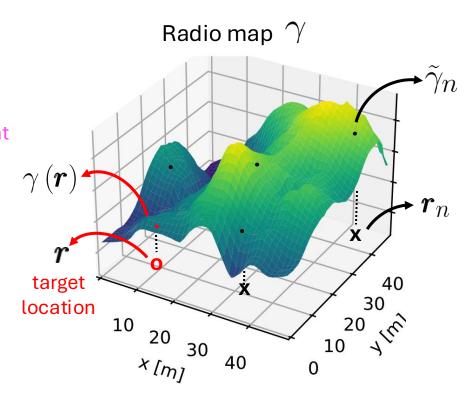
Radio Map Estimation: Problem Formulation

⇔Given



 \circ Target location $~m{r} \in \mathbb{R}^3$





Contribution in RME

Related work

- > Non-DNN
 - KNN
 - Kriging
 - Sparsity-based inference
 - Dictionary learning
 - Kernel-based learning
 - Matrix completion
 - Graphical models
- > DNN
 - U-Net
 - Deep Completion Autoencoders
 - RadioUNet
 - DRUE

[Alaya-Feki et al. 2008] [Bazerque et al. 2010] [Kim et al. 2013] [Romero et al. 2017]

[Schaufele et al. 2019]

[Ha et al. 2024]

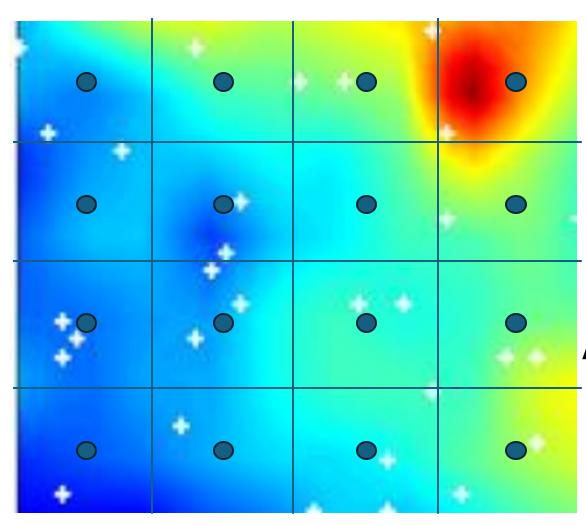
- Oversimplify propagation characteristics
- Cannot learn complex patterns from data

- [Krijestorac et al. 2021] [Teganya et al. 2021] [Levie et al. 2021] [Shrestha et al. 2022]
- Grid discretization
 - No translation equivariance
 - No rotation equivariance
 - Limited spatial resolution
- Large number of parameters

Contributions: Spatial TransfOrmer for Radio Map estimation (STORM)

- 1. Full spatial resolution (gridless) -> translation equivariance + rotation equivariance
- 2. Estimate only at desired locations
- 3. Accommodate measurements outside the region of interest
- 4. Low computational complexity
- 5. Estimation performance sets the state of the art

Gridless vs. grid-based RME



Existing deep learning estimators rely on a grid

Disadvantages of grid discretization

- Limited spatial resolution
- Polynomial growth of the grid size
- Lack of translation equivariance (even discontinuous!)
- Lack of rotation equivariance

Additional limitations of grid-based estimators

- They estimate the map at all grid points
- Cannot accommodate measurements out of the estimation region

Function estimation formulation + transformer \rightarrow Gridless estimator

Attention-Based RME

Spatial TransfOrmer for Radio Map estimation (STORM)

Desiderata

Arbitrary resolution

Invariant to permutations, translation, rotation

Accommodates an arbitrary number of measurements

 $oldsymbol{U}\left(heta^{*}
ight)\left(oldsymbol{r}_{n}
ight)$

 $\cos\left(\theta_{n}\right)$

 $\sin (\theta_n)$

 $\theta^* := \angle \sum \exp(\tilde{\gamma}_n) (\boldsymbol{r}_n - \boldsymbol{r})$

Rotation

(ensures invariance)

Key idea:

transformer for function estimation

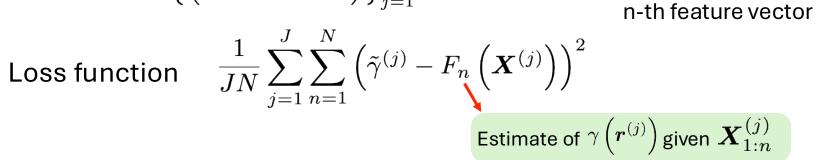
Feature design

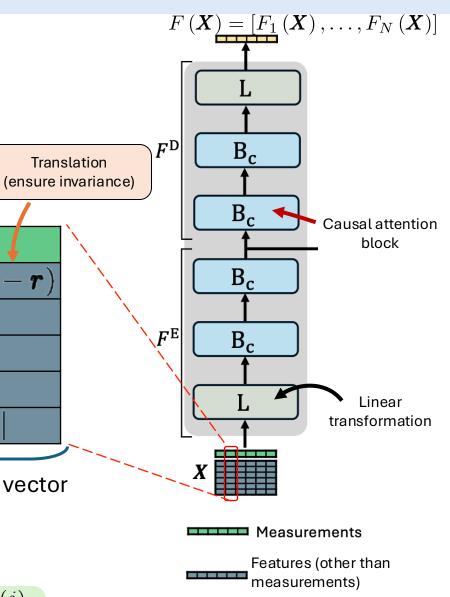
$$\{(\boldsymbol{r}_n, \tilde{\gamma}_n)\}_{n=1}^N$$

Num. realizations

- Training
- Dataset:

$$\left\{\left(oldsymbol{X}^{(j)},oldsymbol{r}^{(j)}, ilde{\gamma}^{(j)}
ight)
ight\}_{j=1}^{J}$$





Estimates

Experimental Setup

STORM:

Parameter	Value
Num. heads	2
Embedding dimension	48
Num. layers	4
~100 k paramters	

❖ Benchmarks*:

Name		Num. Param	
KNN			
Kriging			[Alaya-Feki et al. 2008]
KRR	Kernel ridge regression		[Romero et al. 2017]
DNN-1	Deep Completion Autoencoders	140 k	[Teganya et al. 2020]
DNN-2	RadioUNet	9 M	[Levie et al. 2022]
DNN-3	U-Net	9 M	[Krijestorac et al. 2021]
DNN-4	Autoencoders	60 M	[Shrestha et al. 2022]

^{*} The benchmarks are trained as described in [Shrestha et al. 2024]

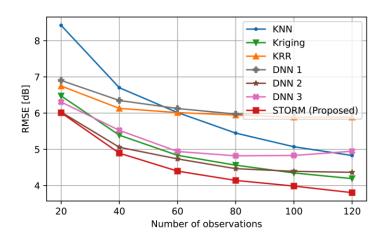
❖ Datasets:

Name	
Ray-tracing	[Teganya and Romero, 2020]
USRP	[Shrestha et al. 2023]
4G Gradiant	[Shrestha et al. 2024]

❖ Performance metric:

$$\text{RMSE} \triangleq \sqrt{\mathbb{E}\left[\frac{1}{|\mathcal{N}_{\text{nobs}}|} \sum_{n \in \mathcal{N}_{\text{nobs}}} \left| \tilde{\gamma}_n - \hat{\gamma}\left(\boldsymbol{r}_n\right) \right|^2\right]}$$
Set of non-observed locations

Experimental Results



→ KNN Kriging 6.0 → DNN 1 RMSE [dB] 8.5 2.7 → DNN 2 DNN 3 DNN 4 STORM (Proposed) 5.6 5.5 5.4 70 80 90 100 50 110 Number of observations

Fig. 3: RMSE for ray-tracing data vs. N when L=64 m, $\Delta=4$ m, and the estimators are trained with $N\in[20,100]$.

Fig. 4: RMSE for USRP data vs. the N when L=38.4 m, $\Delta=1.2$ m, and the estimators are trained with $N\in[40,100]$.

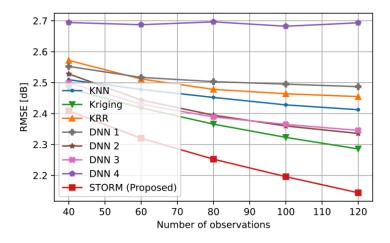


Fig. 5: RMSE for 4G data vs. N when L=64 m, $\Delta=4$ m, and the estimators are trained with $N\in[20,100]$.

Contribution in Active Sensing

Related work

> Non-DNN

- Uncertainty metric + shortest-path algorithm
- Random-active sampling selection
- Uncertainty-aware dynamic programming
- K-means clustering + maximum a posteriori

[Romero et al. 2020]

[Liu et al. 2024]

[Chen et al. 2025]

[Chen et al. 2025]

- Grid discretization
 - No translation equivariance
 - No rotation equivariance
 - Limited spatial resolution
- Ray-tracing simulation dependence, limited applicability

> DNN

- DRUE
- Kriging + Deep Q learning
- Convolutional autoencoder + A*
- Uncertainty-Aware U-Net

[Shrestha et al. 2022]

[Krijestorac et al. 2023]

[Chen et al. 2025]

[Shi et al. 2025]

- Grid discretization
 - No translation equivariance
 - No rotation equivariance
 - Limited spatial resolution
- Large number of parameters

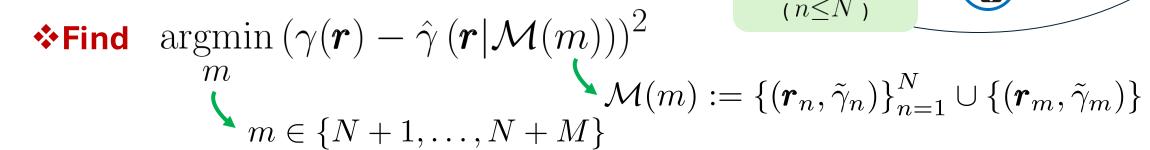
Contributions: Extended transformer-based network for *active sensing*

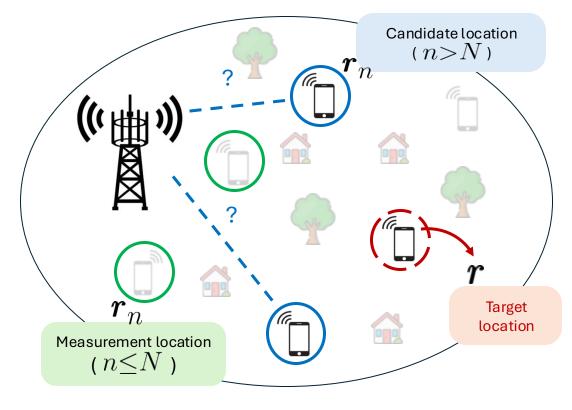
- Attention-based selection of the next measurement location
- Applicable to minimization of drive tests (MDT)
- Low computational complexity

Active Sensing: Problem Formulation

❖Given

- Measurements
- $\{(\boldsymbol{r}_n, \tilde{\gamma}_n)\}_{n=1}^N$
- \circ Candidate locations $\{m{r}_n\}_{n=N+1}^{N+M}$
- \circ Target location $\ r$

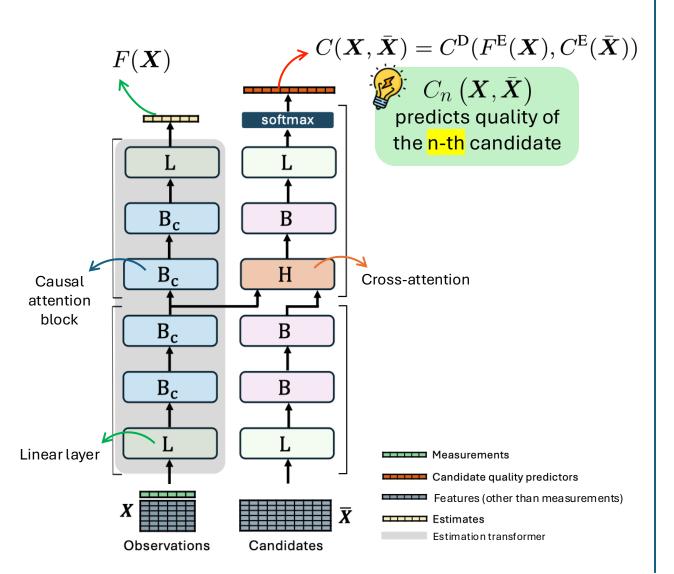


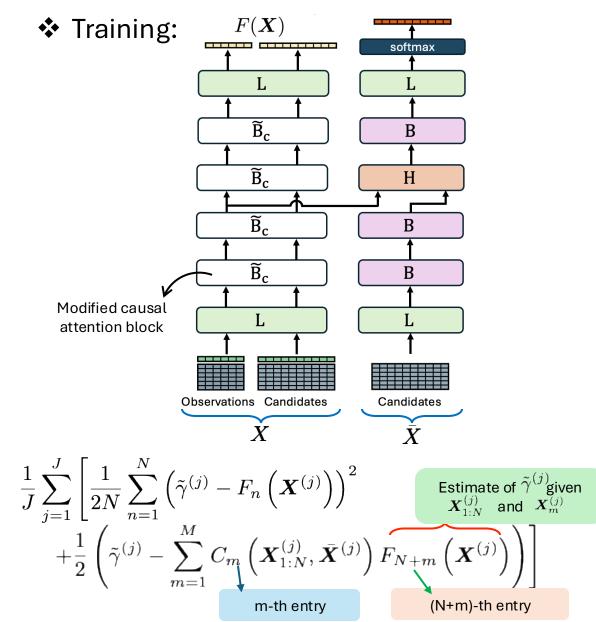


$$=\left\{\left(oldsymbol{r}_{n}, ilde{\gamma}_{n}
ight)
ight\}_{n=1}^{N}\cup\left\{\left(oldsymbol{r}_{m}, ilde{\gamma}_{m}
ight)
ight\}$$

Attention-Based Active Sensing

Inference mode





Experimental Results

STORM:

Parameter	Value
Num. heads	2
Embedding dimension	48
Num. layers	4
~100 k paramters	

❖ Datasets:

Name	
Ray-tracing	
USRP	[Shrestha et al. 2023]

Performance metric:

$$\mathbf{RMSE} \triangleq \sqrt{\mathbb{E}\left[\frac{1}{|\mathcal{N}_{\mathrm{nobs}}|} \sum_{n \in \mathcal{N}_{\mathrm{nobs}}} |\tilde{\gamma}_n - \hat{\gamma}\left(\boldsymbol{r}_n\right)|^2\right]}$$
Set of non-observed locations

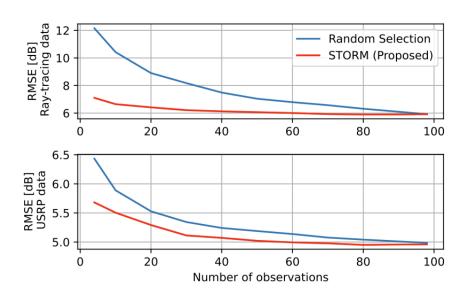
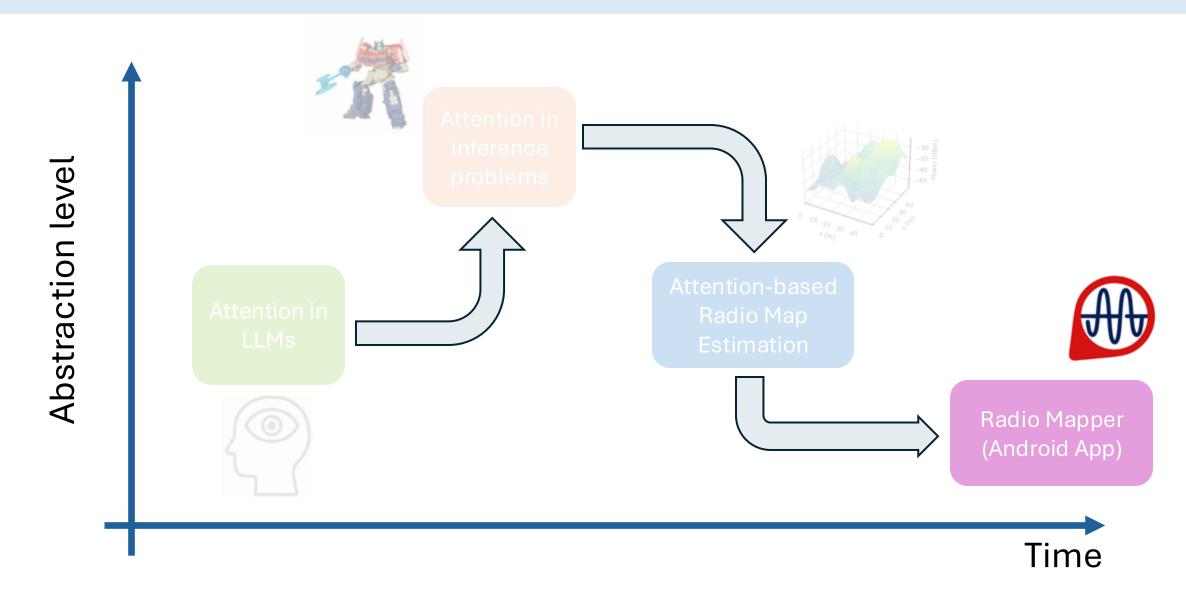


Fig. 6: RMSE vs. N for the active sensing problem with ray-tracing and USRP data. D=20.

Outline of this talk



Radio Mapper

- App for RME with smartphones
 - Measurement collection
 - Map visualization
- Uses STORM (among other estimators)
 - Measurements sent to the server
 - Estimation in the server
- Mapped metrics
 - > 5G -> CSI-RSRP, CSI-RSRQ, CSI-SINR, SS-RSRP, SS-RSRQ, SS-SINR
 - ➤ LTE → RSRP, RSRQ, SINR, CQI
 - ➤ GSM → RSSI



A Mobile App for Radio Map Estimation



Radio Mapper





> Android: Testing phase > download at <u>rme.uia.no</u>





 \rightarrow iOS: No \rightarrow limitations of Apple devices

Challenges

- \rightarrow Earth is round \rightarrow Mercator system \rightarrow Estimation tiles \rightarrow Zoom regions
- ➤ Computation
 Fast response → Multi-threading, code optimization...
 GPU → memory, race conditions...
- Storage
 Database access
- > Security, privacy...

Radio Mapper



Development team





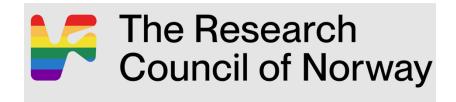


Daniel Romero

Salatiel Genol Paola Quintero

Alejandro Vidal

Sponsors





Conclusions and Future Work

Conclusions:

- Attention is a general inference tool applicable in many different problems
- Spatial TransfOrmer for Radio Map estimation (STORM)
- Active sensing
- Radio Mapper

Future work:

- RME
 - Experiments with indoor radio maps.
 - Accommodate side information, e.g. building locations.
- Active sensing
 - Metrics that account for estimation error on a set of locations
 - > A more extensive experimental study

