# Seminar in Deep Neural Networks 2025

Frédéric Berdoz, Benjamin Estermann
18th of February 2025

# Introduce yourself!

- Name
- Degree, Background in Machine Learning (theoretical and/or practical)
- What are your expectations for the seminar?
- What do you want to learn?

# Supervisors

Andreas

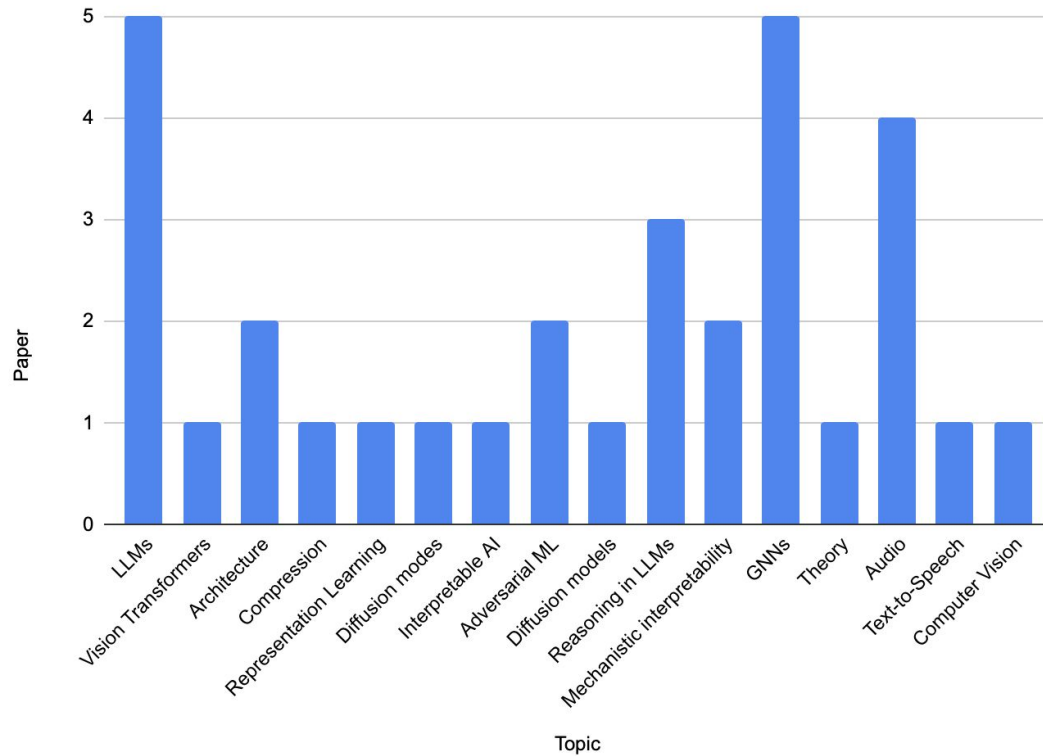Benjamin

Frédéric

Florian

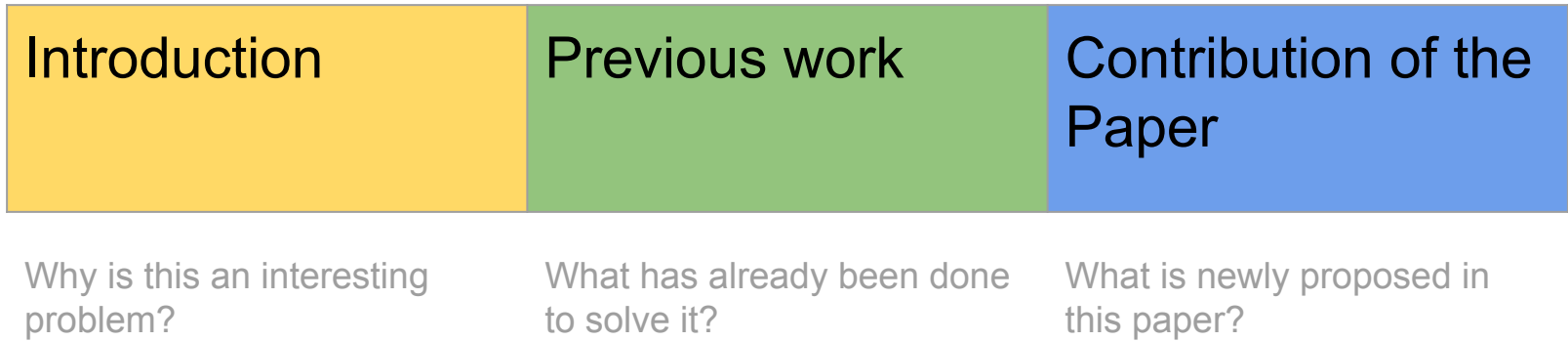Till

Luca

Saku

Sam

# Topic overview

# Schedule

| Date | Presenters | Papers | Mentors | |
|------|-----------|--------|---------|---|
| February 25 | Qi Ma<br>Harald Semmelrock | InstructPix2Pix: Learning To Follow Image Editing Instructions<br>rStar-Math: Small LLMs Can Master Math Reasoning with Self-Evolved Deep Thinking | Florian Grötschla<br>Benjamin Estermann | TBA |
| March 04 | Jakob Hütteneder<br>Yanik Künzi | Guiding a Diffusion Model with a Bad Version of Itself<br>Scaling the Codebook Size of VQGAN to 100,000 with a Utilization Rate of 99% | Till Aczel<br>Luca Lanzendörfer | TBA |
| March 11 | Alexandre Elsig<br>Adam Suma | Towards Compositional Adversarial Robustness: Generalizing Adversarial Training to Composite Semantic Perturbations<br>DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning | Andreas Plesner<br>Samuel Dauncey | TBA |
| March 18 | Niccolò Avogaro<br>Valentin Abadie | Vision Transformers Need Registers<br>Towards Foundation Models for Knowledge Graph Reasoning | Frédéric Berdoz<br>Florian Grötschla | TBA |
| March 25 | Sebastian Brunner<br>Coralie Sage | It's Not What Machines Can Learn, It's What We Cannot Teach<br>Scaling Monosemanticity: Extracting Interpretable Features from Claude 3 Sonnet | Saku Peltonen<br>Samuel Dauncey | TBA |
| April 01 | Anna Kosovskaia<br>Florian Zogaj | You Only Cache Once: Decoder-Decoder Architectures for Language Models<br>Multimodal Neurons in Artificial Neural Networks | Benjamin Estermann<br>Andreas Plesner | TBA |
| April 08 | Diego Arapovic<br>Lukas Rüttgers | Beyond Autoregression: Discrete Diffusion for Complex Reasoning and Planning<br>Convolutional Differentiable Logic Gate Networks | Andreas Plesner<br>Till Aczel | TBA |
| April 15 | Giovanni De Muri<br>Jonas Mirlach | In-context Learning and Induction Heads<br>Good, Cheap, and Fast: Overfitted Image Compression with Wasserstein Distortion | Samuel Dauncey<br>Till Aczel | TBA |
| April 22 | - | Easter Break | - | - |

. . .

ETH zürich

# How to structure your talk

| Introduction | Previous work | Contribution of the Paper |
|:---|:---|:---|
| Why is this an interesting problem? | What has already been done to solve it? | What is newly proposed in this paper? |

# Presentation Style

ETH*zürich*

# How to approach your paper



ETH zürich

# Admin stuff