

# Multilingual Radiology Report Classification

Desmond Elliott

Language and Multimodal Processing Group  
Department of Computer Science  
University of Copenhagen



# Medical

This talk

Everyone else here today

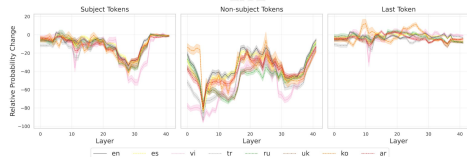
# Multimodal

## How Do Multilingual Language Models Remember Facts?

Constanza Fierro<sup>†</sup> Negar Foroutan<sup>‡</sup> Desmond Elliott<sup>†</sup> Anders Søgaard<sup>†</sup>

<sup>†</sup> Department of Computer Science, University of Copenhagen

<sup>‡</sup> EPFL



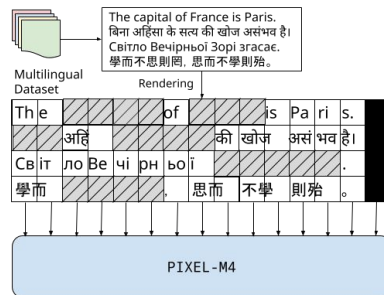
# Multilingual

## Multilingual Pretraining for Pixel Language Models

Ilker Kesen<sup>†</sup> Jonas F. Lotz<sup>†,‡</sup> Ingo Ziegler<sup>†</sup> Phillip Rust<sup>†</sup> Desmond Elliott<sup>†</sup>

<sup>†</sup> Department of Computer Science, University of Copenhagen

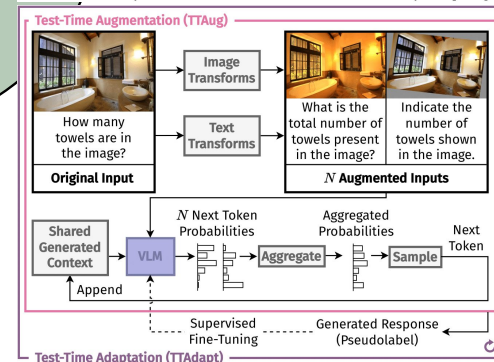
<sup>‡</sup> ROCKWOOL Foundation Research Unit



## EFFICIENT TEST-TIME SCALING FOR SMALL VISION-LANGUAGE MODELS

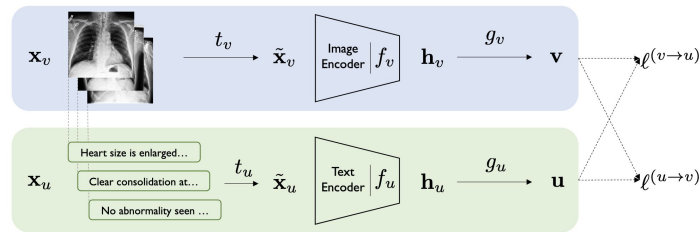
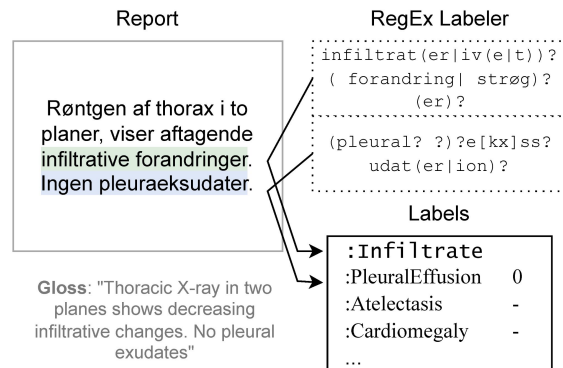
Mehmet Onurcan Kaya<sup>1,2</sup> Desmond Elliott<sup>3,2</sup> Dim P. Papadopoulos<sup>1,2</sup>

<sup>1</sup> Technical University of Denmark <sup>2</sup> Pioneer Center for AI <sup>3</sup> University of Copenhagen



# Radiology Report Classification

- Backbone of training imaging classification systems
  - regex is everywhere
  - SSL emerging as an alternative
    - This is compute-intensive compared to LLM knowledge
- Not much publicly shared data
  - MIMIC-CXR, CheXpert, etc.
- Disjoint findings labels
  - MIMIC-CXR: 15 findings
  - PadChest: 49 findings



# When are radiology reports useful for training medical image classifiers?

Herman Bergström<sup>\*1</sup>, Zhongqi Yue<sup>1</sup>, and Fredrik D. Johansson<sup>1</sup>

<sup>1</sup>Department of Computer Science & Engineering,  
Chalmers University of Technology and University of Gothenburg

---

## MR-CLIP: Efficient Metadata-Guided Learning of MRI Contrast Representations

Mehmet Yigit Avci<sup>1</sup>, Pedro Borges<sup>1</sup>, Paul Wright<sup>1</sup>, Mehmet Yigitsoy<sup>2</sup>,  
Sebastien Ourselin<sup>1</sup>, and Jorge Cardoso<sup>1</sup>

<sup>1</sup> School of Biomedical Engineering and Imaging Sciences, King's College London,  
London, UK

<sup>2</sup> deepc GMBH, Munich, Germany

---

## Are Large Vision Language Models Truly Grounded in Medical Images? Evidence from Italian Clinical Visual Question Answering

---

**Federico Felizzi<sup>1,\*</sup>, Olivia Riccomi<sup>1</sup>, Michele Ferramola<sup>2</sup>, Francesco Andrea Causio<sup>3,1</sup>,  
Manuel Del Medico<sup>3,1,\*</sup>, Vittorio De Vita<sup>3,1</sup>, Lorenzo De Mori<sup>1,4</sup>, Alessandra Piscitelli<sup>1,5</sup>,  
Pietro Eric Risuleo<sup>3,1</sup>, Bianca Destro Castaniti<sup>1,5</sup>, Antonio Cristiano<sup>3,1</sup>,  
Alessia Longo<sup>6</sup>, Luigi De Angelis<sup>1,7</sup>, Mariapia Vassalli<sup>1,5</sup>, Marcello Di Pumpo<sup>3,1</sup>**

<sup>1</sup>SIAM, Rome, Italy   <sup>2</sup>NSBProject, Mantova, Italy

<sup>3</sup>Dept. of Life Sciences & Public Health, UCSC, Rome, Italy

<sup>4</sup>ASL RM 4, Bracciano, Italy   <sup>5</sup>UCSC, Rome, Italy

<sup>6</sup>Univ. Paris Cité, France   <sup>7</sup>Univ. of Pisa, Italy   \*Corresp. author: federico.felizzi@gmail.com

How can we combine  
publicly available radiology  
report resources into a  
single classification model?



Alice

# **MOSAIC: A Multilingual, Taxonomy-Agnostic, and Computationally Efficient Approach for Radiological Report Classification**

**Alice Schiavone<sup>1,2</sup>, Marco Fraccaro<sup>3</sup>, Lea Marie Pehrson<sup>1,4,5</sup>, Silvia Ingala<sup>4,6</sup>, Rasmus Bonnevie<sup>3</sup>  
Michael Bachmann Nielsen<sup>5</sup>, Vincent Beliveau<sup>7</sup>, Melanie Ganz<sup>1,2</sup>, Desmond Elliott<sup>1</sup>**

<sup>1</sup>Department of Computer Science, University of Copenhagen

<sup>2</sup>Neurobiology Research Unit, Copenhagen University Hospital

<sup>3</sup>Unumed Aps, <sup>4</sup>Department of Diagnostic Radiology, Copenhagen University Hospital

<sup>5</sup>Department of Clinical Medicine, University of Copenhagen

<sup>6</sup>Cerebriu A/S, <sup>7</sup>Institute for Human Genetics, Medical University of Innsbruck

# Desiderata

---

- **Fully open source:** keep your medical data on-site
- **Accessible:** run and train on inexpensive general-purpose GPUs
  - training and inference on an 24GB RTX 3090
- **Multilingual:** works for any EU26 major language
  - only evaluated in 4 languages due to data availability
- **Flexible:** adapts to different findings labels with minimal intervention
  - LLMs are reservoirs of written human knowledge

# The LLM money pit

The rich man experience:

- LLM can solve 6 additional **high-school** competition math problems (AIME) for 4.2M USD
  - Reaching 65% on the test (below the acceptance cutoff) using only 16 tries...

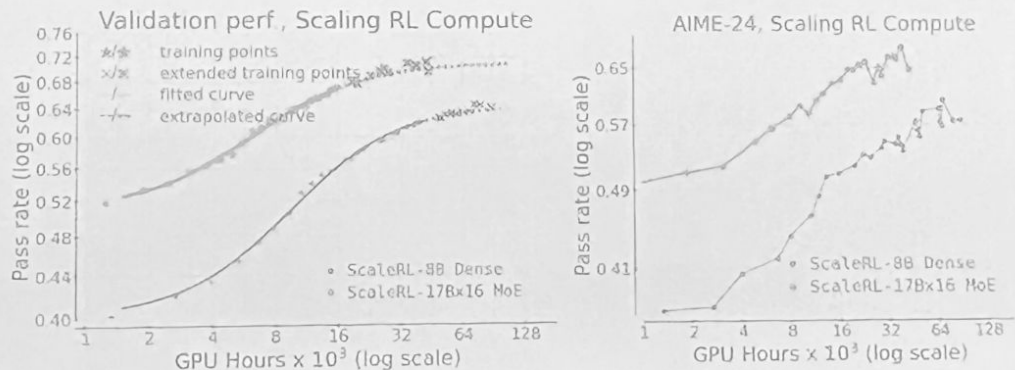


Figure 1 Predictably Scaling RL compute to 100,000 GPU Hours (a) We ran ScaleRL for 100k GPU hours on an 8B dense

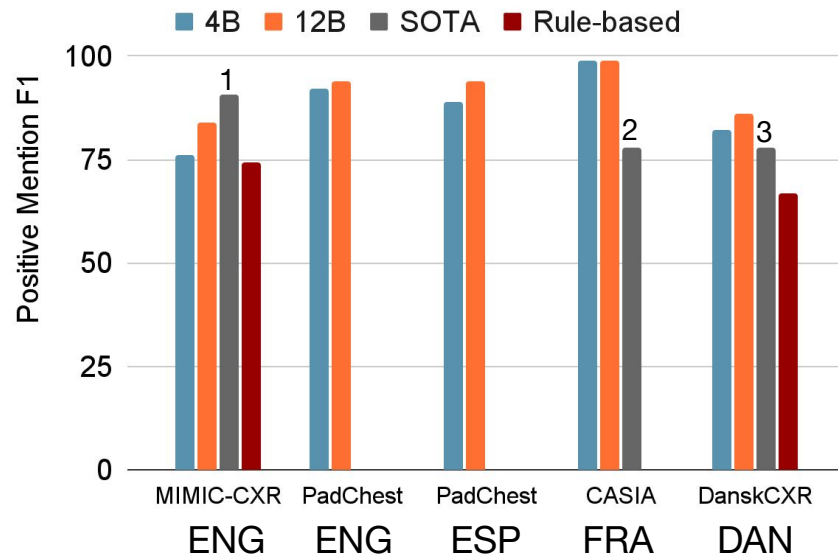
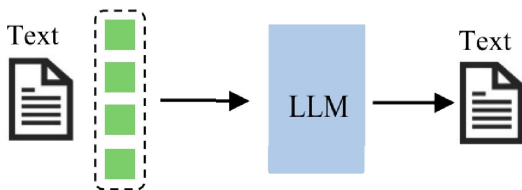
Alexia Jolicoeur-Martineau, Mila, October 2025





# MOSAIC-4B and 12B

- Finetuned on 10K reports in English, Spanish, and French
  - QLoRA optimization on the Q,K,V, FF, and Output layers
  - Need maximum of 16.2G VRAM and 33 minutes for SFT
- Prompt-based inference that can predict up to 68 findings



1. CheX-GPT, 2. CASIA-CLS, 3. DanskBERT

*Lucas Dixon, Google DeepMind*  
How <sup>↑</sup> I think  
about LLMs...

**An interpreter** (that can translate  
between languages, concepts, and styles)



**An improv  
comedian**

**A fuzzy  
database of  
the web**

# Prompt-based Inference

---

Require JSON-structured responses

You are a helpful radiology assistant. Given a radiology report, classify each abnormality into a class. Output a valid JSON with each abnormality as key, and the class as value. The keys must be {findings}. The values can be one of {classes}. The values have the following interpretation:

Define style of positive/uncertain findings

(1) the abnormality was mentioned, even with uncertainty, in the report, e.g. 'A large pleural effusion', 'The cardiac contours are stable.', 'The cardiac size cannot be evaluated.';

Negative mentions

(2) the abnormality was negatively mentioned in the report; e.g. 'No pneumothorax.'

# Datasets

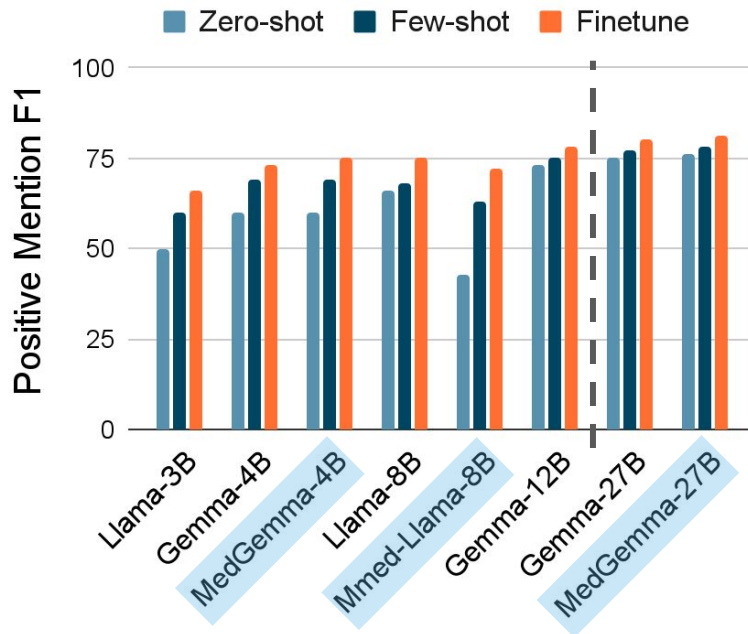
---

Dataset	Language	Modality	Number of Findings	Avg. Chars	Mention Classes	Train	Dev	Test
MIMIC-CXR	en	Chest X-Ray	14	760	+, -, ~	535	50	100
PadChest-GR	es, en	Chest X-Ray	49	115	+	1951	100	879
CASIA-CXR	fr	Chest X-Ray	5	400	+	7677	100	3334
DanskCXR	da	Chest X-Ray	48	312	+, -	1600	125	750
DanskMRI	da	Brain MRI	3	1941	+, -, ~	194	50	345

- Focus on publicly available datasets
  - 194-7600 training examples
  - 115–1941 characters
  - 3–49 findings across variable number of mention classes
- DanskMRI evaluates performance on different imaging modality

# Which Backbone LLM?

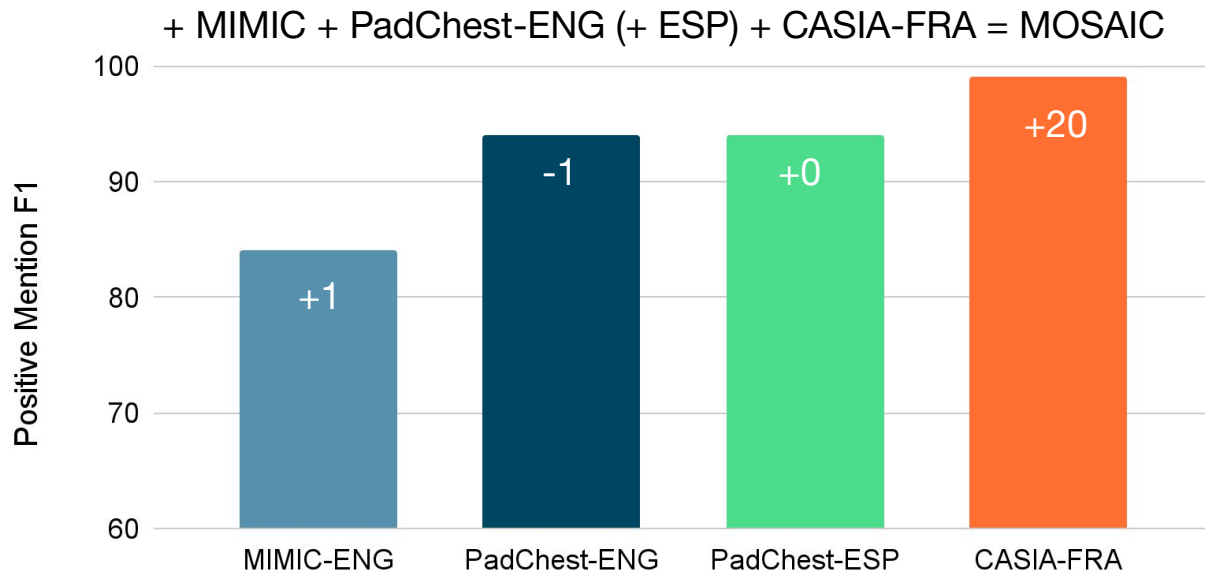
- Setups:
  - Zero-shot prompting
  - Few-shot prompting
  - Dataset-specific fine-tuning
- Gemma and LLaMA LLMs
  - 3B–27B variants
  - General and **medical domain**



Finding 1: No substantial difference between  
general / medical domain models

# Finetuning on Public Datasets

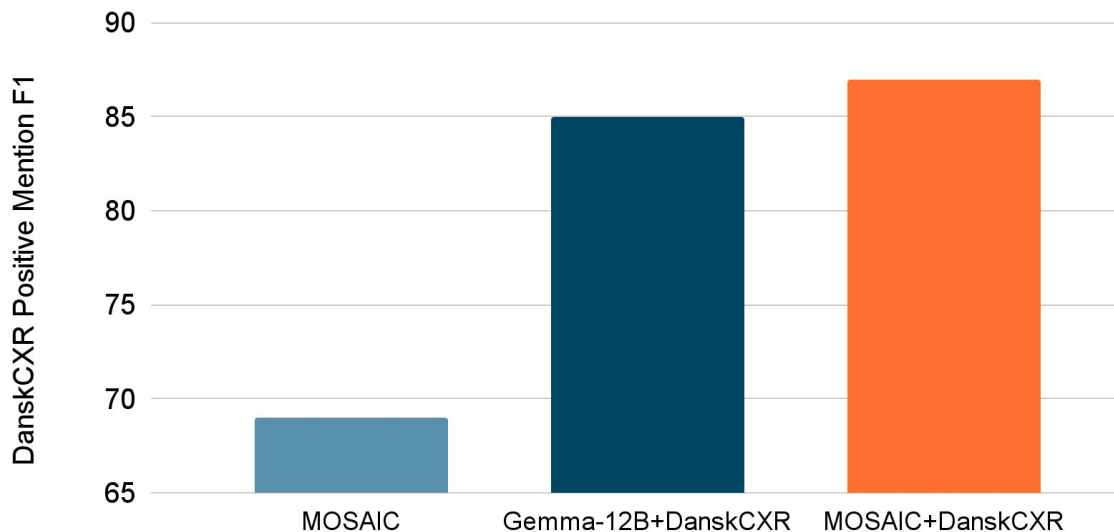
- How does performance improve as we train on different label sets?



Finding 2: Improvements are additive and do not seem to interfere

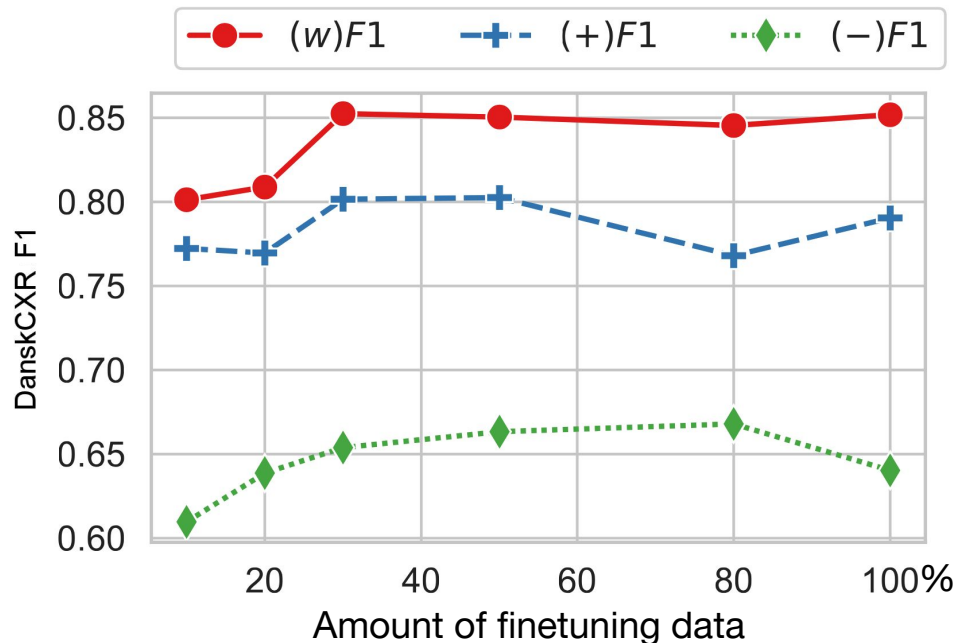
# New Dataset Adaptation

---



Finding 3: MOSAIC is a better starting point for new data

# How Much Data Do You Need?

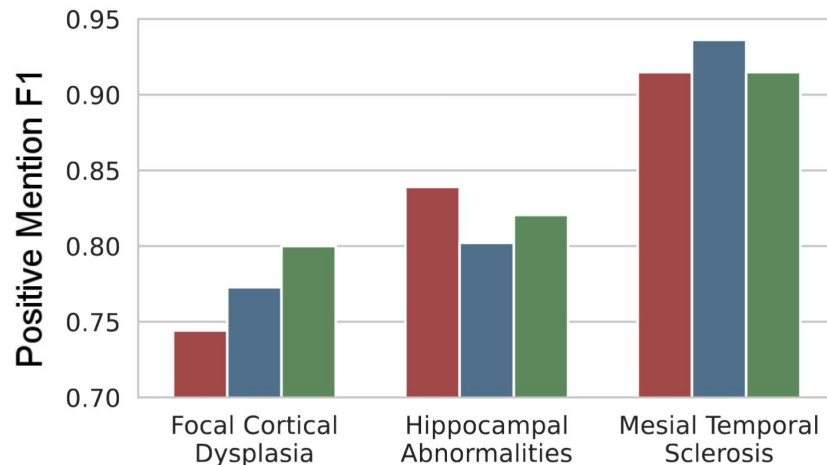
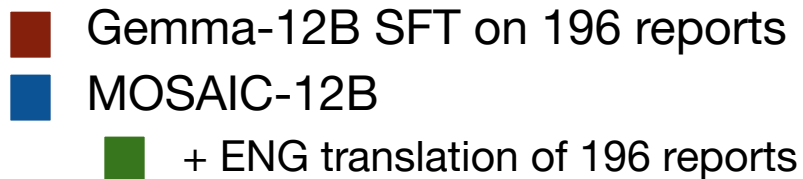


Finding 4: You don't need much data if you start from MOSAIC



# Different Imaging Modality

- Adapt MOSAIC to predicting three findings in Epilepsy MRI reports



Finding 5: MOSAIC can be repurposed to a new modality

# Open Directions

---

- **Multimodal inputs** could improve performance but how to handle reports from different imaging modalities
- **Simple text-only augmentation** could substantially improve performance [Aepli and Sennrich, 2022; Kaya et al. 2025]
- **Multi-agent LLMs** could better handle different mention classes
- **Broken tokenizers** could be fixed to further improve performance
  - See, e.g. TokenDist [Dobler et al. 2025]
- **Synthetic data generation** using self-consistency [Wang et al. 2023]

# Conclusions

---

- Multilingual LLMs are radiology report classifiers
  - Handle different label sets
  - Handle reports from different imaging modalities
- Multilingual multi dataset SFT can reduce the total amount of data that needs expert annotation
  - Focus the time of our clinical colleagues on labelling lower-frequency findings or difficult examples
- MOSAIC is open source
  - Please tell us if it works for your data and language

# References

---

- A. Schiavone M. Fraccaro, L. M. Pehrson, S. Ingala, R. Bonnevie, M. B. Nielsen, V. Beliveau, M. Ganz, and D. Elliott. 2025. MOSAIC: A Multilingual, Taxonomy-Agnostic, and Computationally Efficient Approach for Radiological Report Classification
- K. Dobler, D. Elliott, and G. de Melo. 2025. Token Distillation: Attention-aware Input Embeddings For New Tokens.
- M. O. Kaya, D. Elliott, and D. P. Papadopoulos. 2025. Efficient Test-Time Scaling for Small Vision-Language Models.
- Aepli and Sennrich. ACL 2022. Improving Zero-Shot Cross-lingual Transfer Between Closely Related Languages by Injecting Character-Level Noise.
- Wang et al. ICLR 2023. Self-Consistency Improves Chain of Thought Reasoning in Language Models.