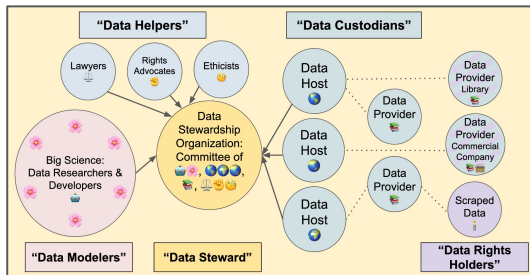# Building Better Language Models: *Insights from BigScience* 🌸

Colin Raffel
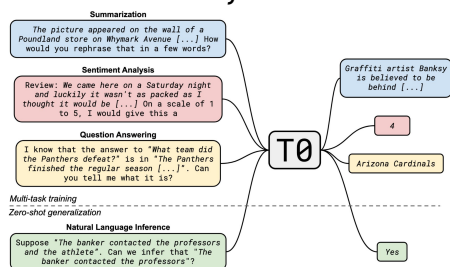
# Data governance
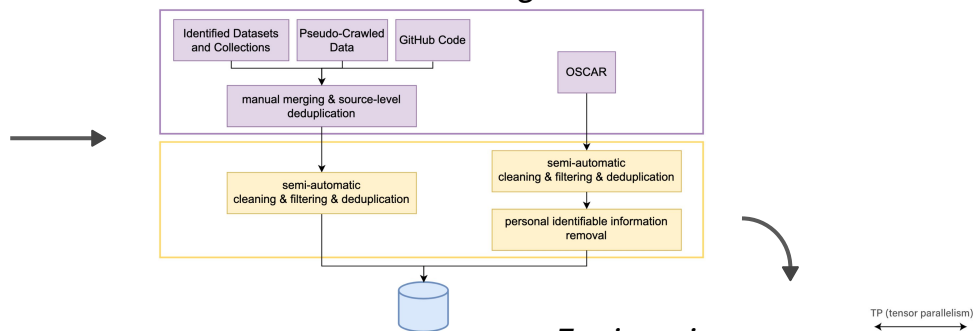


# Data sourcing



# Evaluation



# Engineering



# Objective



# Modeling

From *https://huggingface.co/spaces/bigscience/SourcingCatalog*

8 copies of the model are trained in parallel
on a total of 384 GPUs (data parallelism = 8)

TP (tensor parallelism)

DP (data parallelism)

PP (pipeline parallelism)

Model parameters
are divided across 4 GPUs
(tensor parallelism = 4)

The layers of the model
are spread across
12 groups of GPUs
(pipeline parallelism = 12)

One full copy («replica»)
of the model takes
48 GPUs

data batch #1    data batch #2    data batch #3    data batch #4    data batch #5    data batch #6    data batch #7    data batch #8

data batch

⬜ → 1 GPU - NVIDIA A100 with 80GB of memory

*From "BLOOM: A 176B-Parameter Open-Access Multilingual Language Model" by Le Scao et al.*

## Accuracy ↑

InstructGPT davinci v2 (175B*)
TNLG v2 (530B)
Anthropic-LM v4-s3 (52B)
OPT (175B)
Cohere xlarge v20220609 (52.4B)
J1-Jumbo v1 (178B)
GPT-3 davinci v1 (175B)
GLM (130B)
OPT (66B)
BLOOM (176B)
J1-Grande v1 (17B)
Cohere large v20220720 (13.1B)
GPT-NeoX (20B)
J1-Large v1 (7.5B)
InstructGPT curie v1 (6.7B*)
TNLG v2 (6.7B)
GPT-3 curie v1 (6.7B)
GPT-J (6B)
Cohere medium v20220720 (6.1B)
InstructGPT babbage v1 (1.3B*)
UL2 (20B)
T0pp (11B)
T5 (11B)
YaLM (100B)
GPT-3 babbage v1 (1.3B)
GPT-3 ada v20220720 (410M)
InstructGPT ada v1 (350M*)
Cohere small v20220720 (410M)

## Calibration error ↓

InstructGPT ada v1 (350M*)
OPT (66B)
InstructGPT babbage v1 (1.3B)
InstructGPT curie v1 (6.7B*)
BLOOM (176B)
OPT (175B)
YaLM (100B)
GPT-NeoX (20B)
InstructGPT davinci v2 (175B*)
GPT-J (6B)
UL2 (20B)
T5 (11B)
Cohere small v20220720 (410M)
GPT-3 babbage v1 (1.3B)
Cohere medium v20220720 (6.1B)
Cohere xlarge v20220609 (52.4B)
GPT-3 curie v1 (6.7B)
GPT-3 ada v1 (350M)
GPT-3 davinci v1 (175B)
TNLG v2 (530B)
TNLG v2 (6.7B)
Cohere large v20220720 (13.1B)
GLM (130B)
T0pp (11B)

## Robustness ↑

InstructGPT davinci v2 (175B*)
Anthropic-LM v4-s3 (52B)
GLM (130B)
TNLG v2 (530B)
BLOOM (176B)
OPT (175B)
Cohere xlarge v20220609 (52.4B)
J1-Jumbo v1 (178B)
GPT-3 davinci v1 (175B)
OPT (66B)
J1-Grande v1 (17B)
GPT-NeoX (20B)
J1-Large v1 (7.5B)
Cohere large v20220720 (13.1B)
InstructGPT curie v1 (6.7B*)
UL2 (20B)
GPT-3 curie v1 (6.7B)
GPT-J (6B)
TNLG v2 (6.7B)
Cohere medium v20220720 (6.1B)
T0pp (11B)
InstructGPT babbage v1 (1.3B*)
T5 (11B)
YaLM (100B)
GPT-3 babbage v1 (1.3B)
Cohere small v20220720 (410M)
InstructGPT ada v1 (350M*)
GPT-3 ada v1 (350M)

## Fairness ↑

InstructGPT davinci v2 (175B*)
TNLG v2 (530B)
Anthropic-LM v4-s3 (52B)
OPT (175B)
BLOOM (176B)
Cohere xlarge v20220609 (52.4B)
OPT (66B)
J1-Jumbo v1 (178B)
GPT-3 davinci v1 (175B)
GLM (130B)
J1-Grande v1 (17B)
Cohere large v20220720 (13.1B)
GPT-NeoX (20B)
J1-Large v1 (7.5B)
InstructGPT curie v1 (6.7B*)
TNLG v2 (6.7B)
GPT-J (6B)
GPT-3 curie v1 (6.7B)
Cohere medium v20220720 (6.1B)
InstructGPT babbage v1 (1.3B*)
UL2 (20B)
T0pp (11B)
T5 (11B)
GPT-3 babbage v1 (1.3B)
YaLM (100B)
GPT-3 ada v1 (350M)
Cohere small v20220720 (410M)
InstructGPT ada v1 (350M*)

## Bias ↓

GPT-J (6B)
InstructGPT ada v1 (350M*)
YaLM (100B)
GPT-3 davinci v1 (175B)
InstructGPT curie v1 (6.7B*)
GPT-3 curie v1 (6.7B)
TNLG v2 (6.7B)
GPT-3 babbage v1 (1.3B)
GPT-NeoX (20B)
T5 (11B)
Cohere medium v20220720 (6.1B)
GLM (130B)
Cohere small v20220720 (410M)
Cohere large v20220720 (13.1B)
GPT-3 ada v1 (350M)
J1-Grande v1 (17B)
BLOOM (176B)
UL2 (20B)
Anthropic-LM v4-s3 (52B)
InstructGPT davinci v2 (175B*)
Cohere xlarge v20220609 (52.4B)
TNLG v2 (530B)
InstructGPT babbage v1 (1.3B*)
J1-Large v1 (7.5B)
OPT (175B)
T0pp (11B)
OPT (66B)
J1-Jumbo v1 (178B)

## Toxicity ↓

T0pp (11B)
GPT-J (6B)
J1-Large v1 (7.5B)
Cohere small v20220720 (410M)
Cohere large v20220720 (13.1B)
YaLM (100B)
J1-Grande v1 (17B)
BLOOM (176B)
Anthropic-LM v4-s3 (52B)
GPT-NeoX (20B)
UL2 (20B)
Cohere xlarge v20220609 (52.4B)
Cohere medium v20220720 (6.1B)
GPT-3 curie v1 (6.7B)
GPT-3 ada v1 (350M)
T5 (11B)
InstructGPT davinci v2 (175B*)
GLM (130B)
TNLG v2 (530B)
OPT (66B)
GPT-3 davinci v1 (175B)
GPT-3 babbage v1 (1.3B)
TNLG v2 (6.7B)
InstructGPT babbage v1 (1.3B*)
InstructGPT curie v1 (6.7B*)
J1-Jumbo v1 (178B)
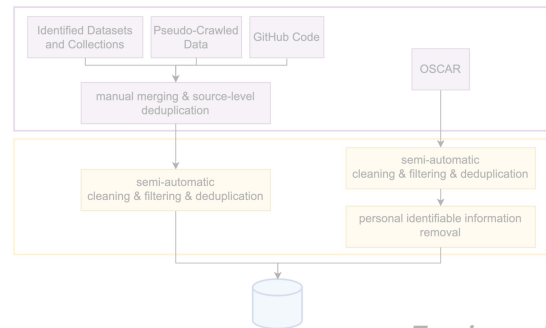InstructGPT ada v1 (350M*)

*From "Holistic Evaluation of Language Models" by Liang et al.*

# Data governance

"Data Helpers"
"Data Custodians"

Lawyers
Rights Advocates
Data Stewardship Organization: Committee of

Data Researchers & Developers

Data Host
Data Provider
Data Provider Library

Data Host
Data Provider
Data Provider Commercial Company

Data Host
Scraped Data

"Data Modelers"
"Data Steward"
"Data Rights Holders"

# Data sourcing

Identified Datasets and Collections | Pseudo-Crawled Data | GitHub Code

OSCAR

manual merging & source-level deduplication

semi-automatic cleaning & filtering & deduplication

semi-automatic cleaning & filtering & deduplication

personal identifiable information removal

# Evaluation

SuperGLUE 0-shot

Ax-b  Ax-b  BoolQ  CB  WiC  WiC

SuperGLUE 1-shot

Ax-b  Ax-b  BoolQ  CB  WiC  WiC

# Engineering

TP (tensor parallelism)

DP (data parallelism)

PP (pipeline parallelism)

data batch #1  data batch #2  data batch #3  data batch #4  data batch #5  data batch #6  data batch #7  data batch #8

☐ ⟶ 1 GPU - NVIDIA A100 with 80GB of memory

# Objective

**Summarization**
*The picture appeared on the wall of a Poundland store on Whymark Avenue [...] How would you rephrase that in a few words?*

**Sentiment Analysis**
*Review: We came here on a Saturday night and luckily it wasn't as packed as I thought it would be [...] On a scale of 1 to 5, I would give this a*

**Question Answering**
*I know that the answer to "What team did the Panthers defeat?" is in "The Panthers finished the regular season [...]". Can you tell me what it is?*

Multi-task training
Zero-shot generalization

**Natural Language Inference**
*Suppose "The banker contacted the professors and the athlete". Can we infer that "The banker contacted the professors"?*

T0

*Graffiti artist Banksy is believed to be behind [...]*

4

*Arizona Cardinals*

*Yes*

# Modeling

Head fusion

Weighted sum of values

Softmax

ALiBi mask

0
-1  0
-2  -1  0

Key-query product

Softmax  Softmax  Softmax

Embed^T  Embed^T  Embed^T

MLP  MLP  MLP

LN  LN  LN

Decoder Block (x 70)

Multi-Head Attention

LN  LN  LN

Embed  Embed  Embed

Token₁  Token₂  Token₃
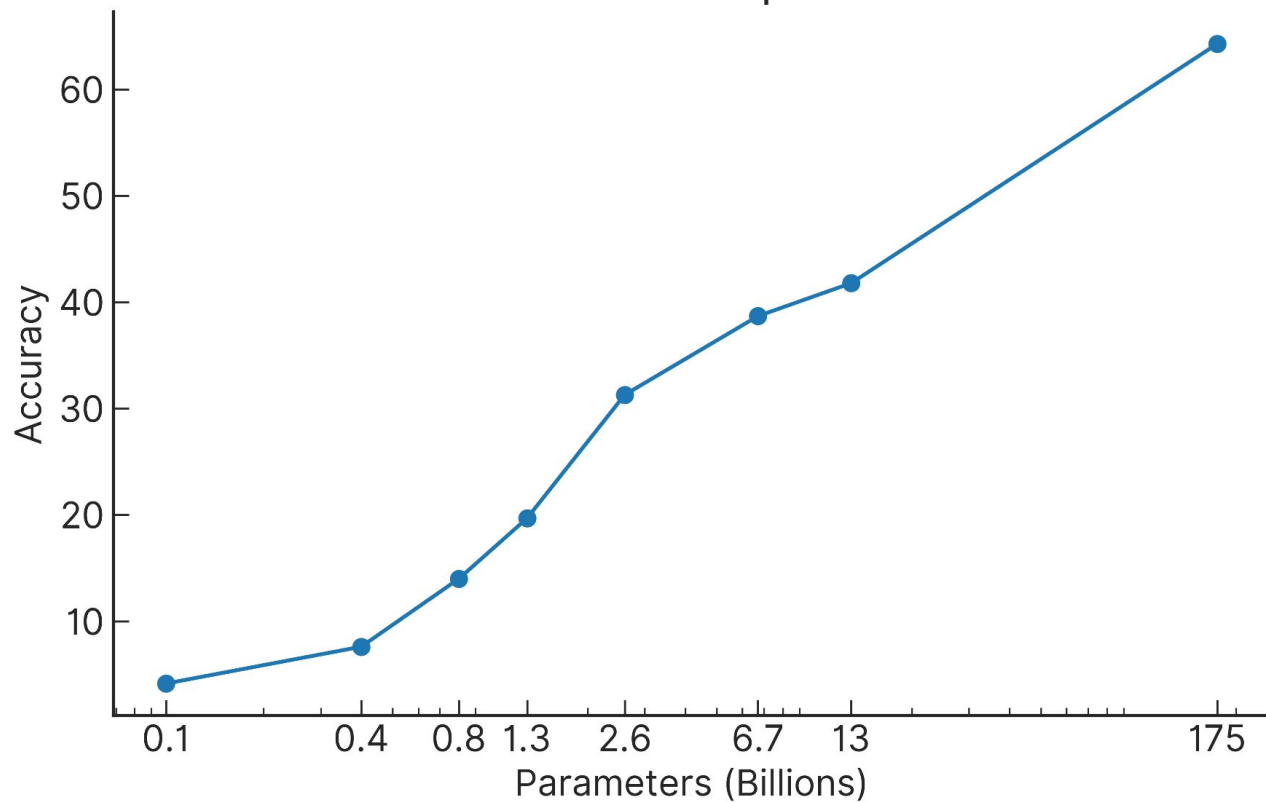
# TriviaQA zero-shot performance



*From "Language Models are Few-Shot Learners" by Brown et al.*

**Closed-book question answering**

http://www.autosweblog.com/cat/trivia-questions-from-the-50s

who was frank sinatra? a: an american singer, actor, and producer.

**Paraphrase identification**

https://www.usingenglish.com/forum/threads/60200-Do-these-sentences-mean-the-same

Do these sentences mean the same? No other boy in this class is as smart as the boy. No other boy is as smart as the boy in this class.

**Natural Language Inference**

https://ell.stackexchange.com/questions/121446/what-does-this-sentence-imply

If I say: He has worked there for 3 years. does this imply that he is still working at the moment of speaking?
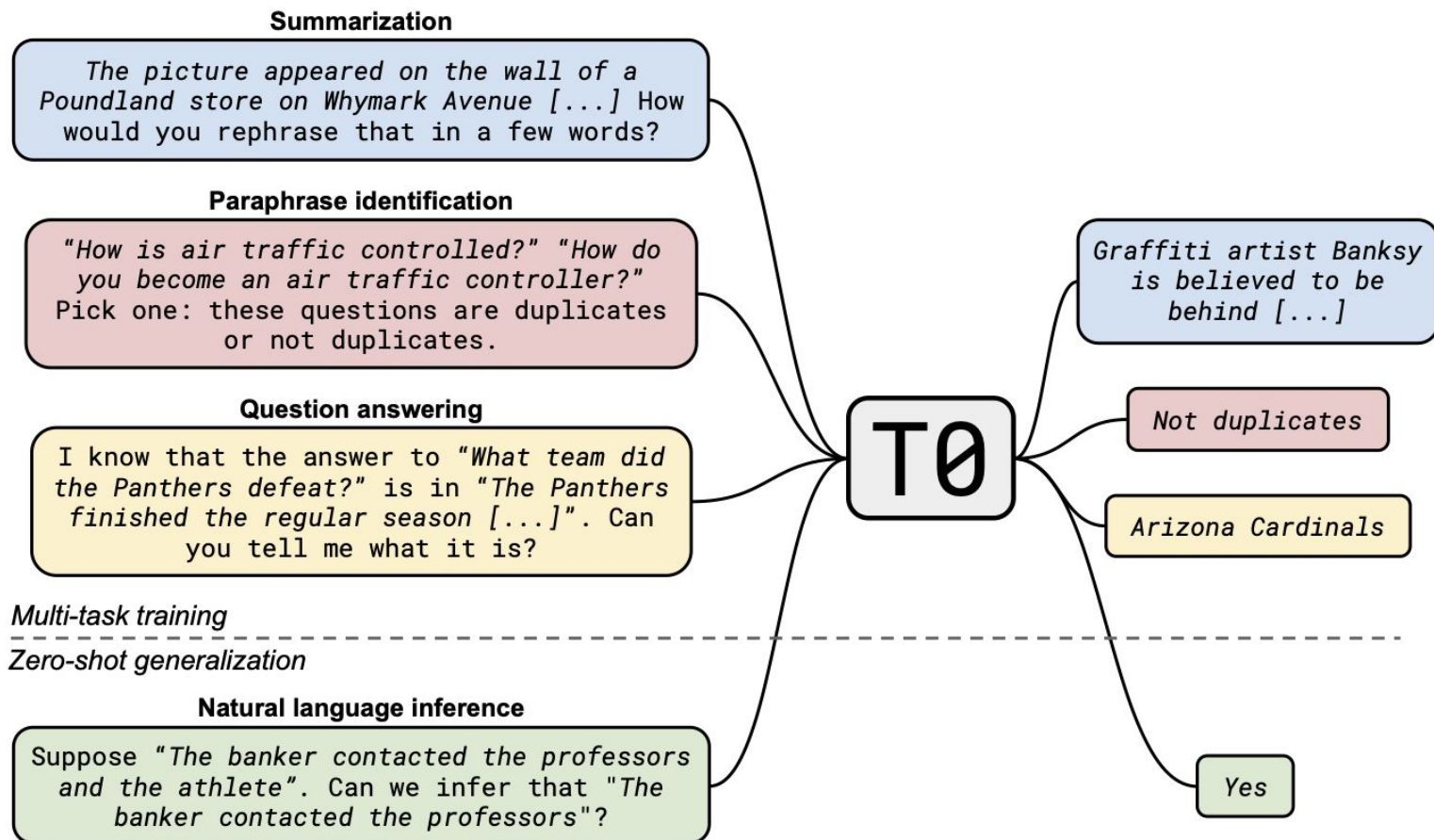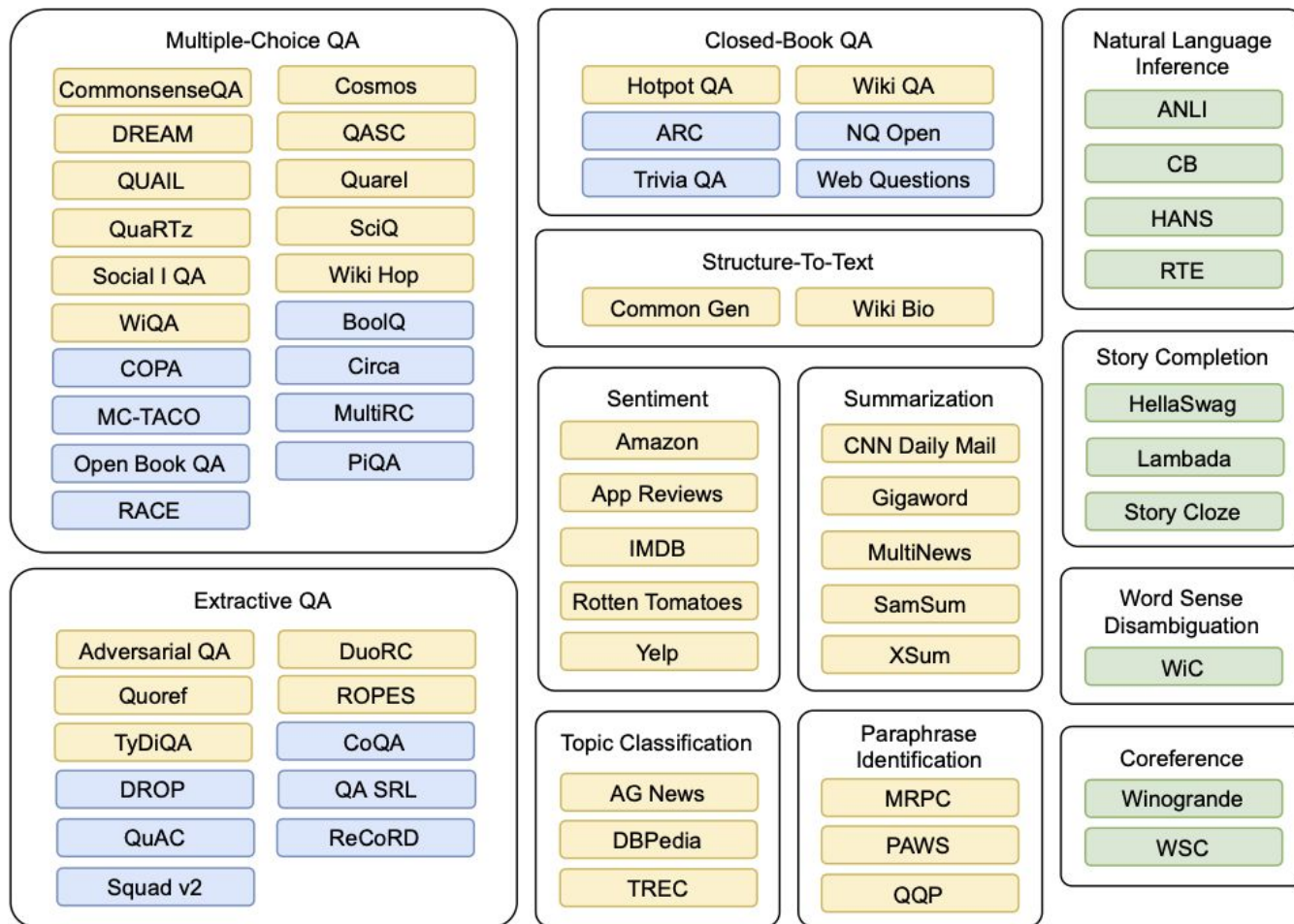
**Summarization**

https://blog.nytsoi.net/tag/reddit

… Lately I've been seeing a pattern regarding videos stolen from other YouTube channels, reuploaded and monetized with ads. These videos are then mass posted on Reddit by bots masquerading as real users. tl;dr: Spambots are posting links to stolen videos on Reddit, copying comments from others to masquerade as legitimate users.

**Pronoun resolution**

https://nursecheung.com/ati-teas-guide-to-english-language-usage-understanding-pronouns/

Jennifer is a vegetarian, so she will order a nonmeat entrée. In this example, the pronoun she is used to refer to Jennifer.

**Summarization**

*The picture appeared on the wall of a Poundland store on Whymark Avenue [...]* How would you rephrase that in a few words?

**Paraphrase identification**

*"How is air traffic controlled?" "How do you become an air traffic controller?"* Pick one: these questions are duplicates or not duplicates.

**Question answering**

I know that the answer to *"What team did the Panthers defeat?"* is in *"The Panthers finished the regular season [...]"*. Can you tell me what it is?

*Multi-task training*
- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
*Zero-shot generalization*

**Natural language inference**

Suppose *"The banker contacted the professors and the athlete"*. Can we infer that *"The banker contacted the professors"*?

**T0**

*Graffiti artist Banksy is believed to be behind [...]*

*Not duplicates*

*Arizona Cardinals*

*Yes*

*From "Multitask Prompted Training Enables Zero-Shot Task Generalization" by Sanh et al.*

## Multiple-Choice QA

| | |
|---|---|
| CommonsenseQA | Cosmos |
| DREAM | QASC |
| QUAIL | Quarel |
| QuaRTz | SciQ |
| Social I QA | Wiki Hop |
| WiQA | BoolQ |
| COPA | Circa |
| MC-TACO | MultiRC |
| Open Book QA | PiQA |
| RACE | |

## Extractive QA

| | |
|---|---|
| Adversarial QA | DuoRC |
| Quoref | ROPES |
| TyDiQA | CoQA |
| DROP | QA SRL |
| QuAC | ReCoRD |
| Squad v2 | |

## Closed-Book QA

| | |
|---|---|
| Hotpot QA | Wiki QA |
| ARC | NQ Open |
| Trivia QA | Web Questions |

## Structure-To-Text

| | |
|---|---|
| Common Gen | Wiki Bio |

## Sentiment

- Amazon
- App Reviews
- IMDB
- Rotten Tomatoes
- Yelp

## Summarization

- CNN Daily Mail
- Gigaword
- MultiNews
- SamSum
- XSum

## Topic Classification

- AG News
- DBPedia
- TREC

## Paraphrase Identification

- MRPC
- PAWS
- QQP

## Natural Language Inference

- ANLI
- CB
- HANS
- RTE

## Story Completion

- HellaSwag
- Lambada
- Story Cloze

## Word Sense Disambiguation

- WiC

## Coreference

- Winogrande
- WSC

*From "Multitask Prompted Training Enables Zero-Shot Task Generalization" by Sanh et al.*

## QQP (Paraphrase)

| Question1 | How is air traffic controlled? |
|-----------|-------------------------------|
| Question2 | How do you become an air traffic controller? |
| Label | 0 |

{Question1} {Question2}
Pick one: These questions
are duplicates or not
duplicates.

→ {Choices[label]}

I received the questions
"{Question1}" and
"{Question2}". Are they
duplicates?

→ {Choices[label]}

## XSum (Summary)

| Document | The picture appeared on the wall of a Poundland store on Whymark Avenue... |
|----------|---------------------------------------------------------------------------|
| Summary | Graffiti artist Banksy is believed to be behind... |

{Document}
How would you
rephrase that in
a few words?

→ {Summary}

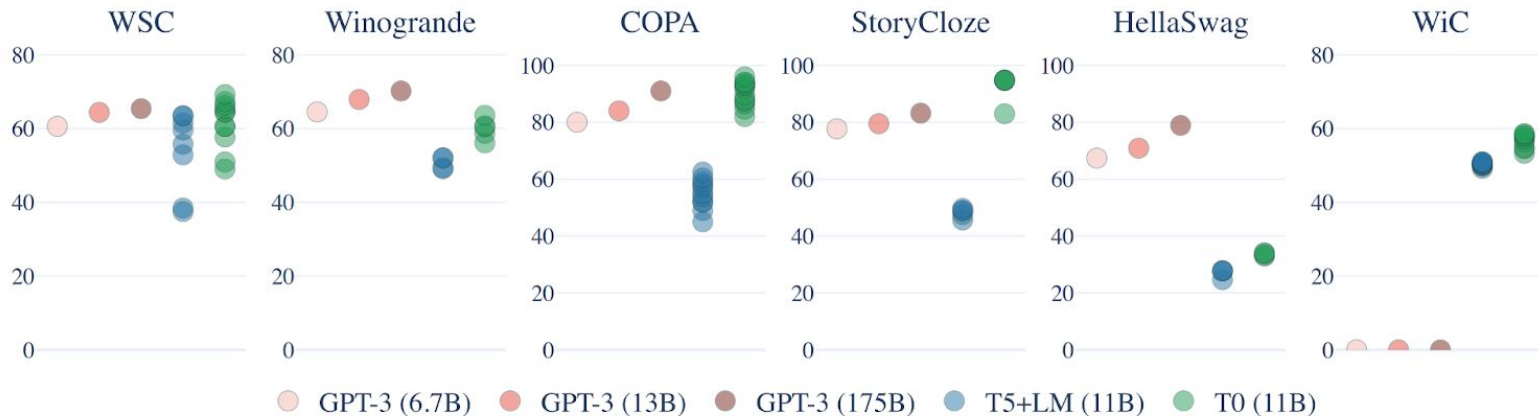First, please read the article:
{Document}
Now, can you write me an
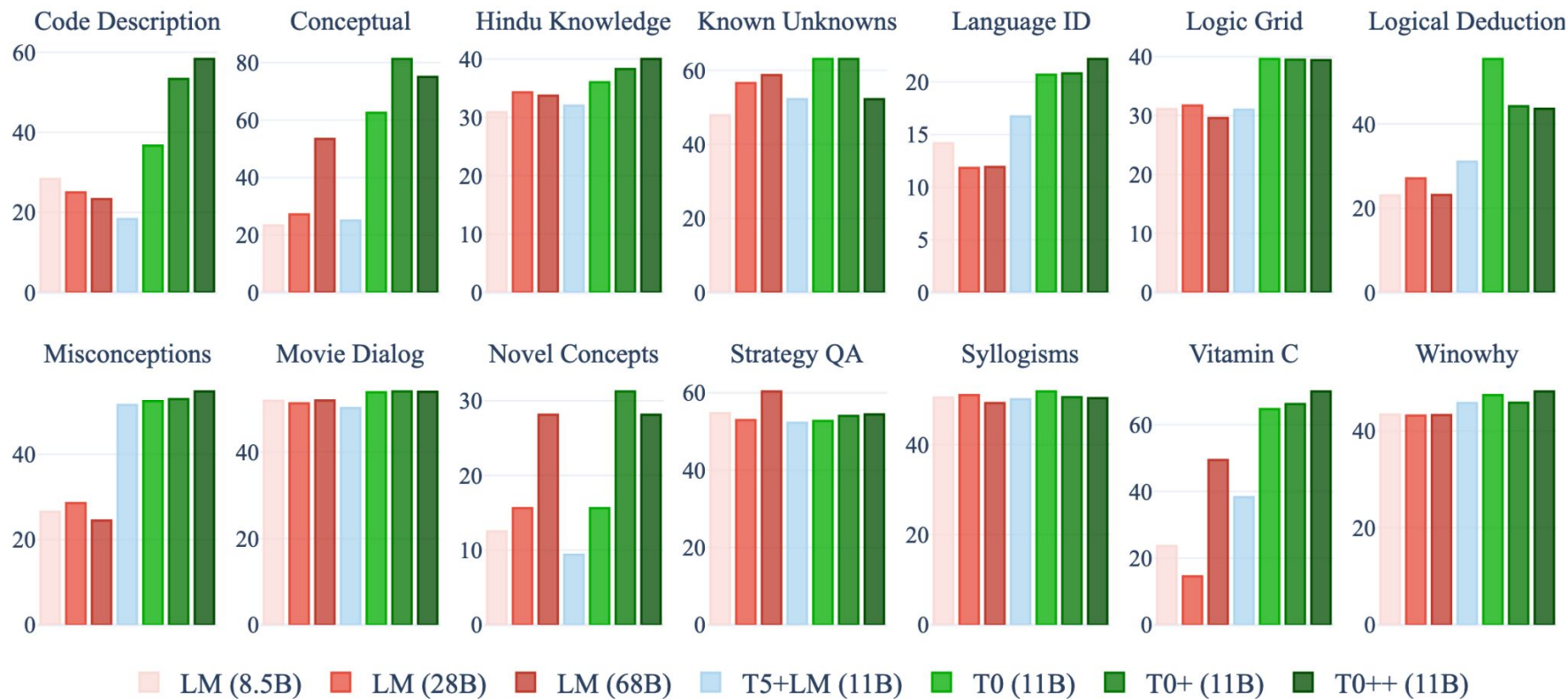extremely short abstract for it?

→ {Summary}

*From "Multitask Prompted Training Enables Zero-Shot Task Generalization" by Sanh et al.*

**Natural Language Inference**

RTE   CB   ANLI R1   ANLI R2   ANLI R3

**Coreference Resolution**

WSC   Winogrande

**Sentence Completion**

COPA   StoryCloze   HellaSwag

**Word Sense**

WiC

○ GPT-3 (6.7B)   ● GPT-3 (13B)   ● GPT-3 (175B)   ● T5+LM (11B)   ● T0 (11B)

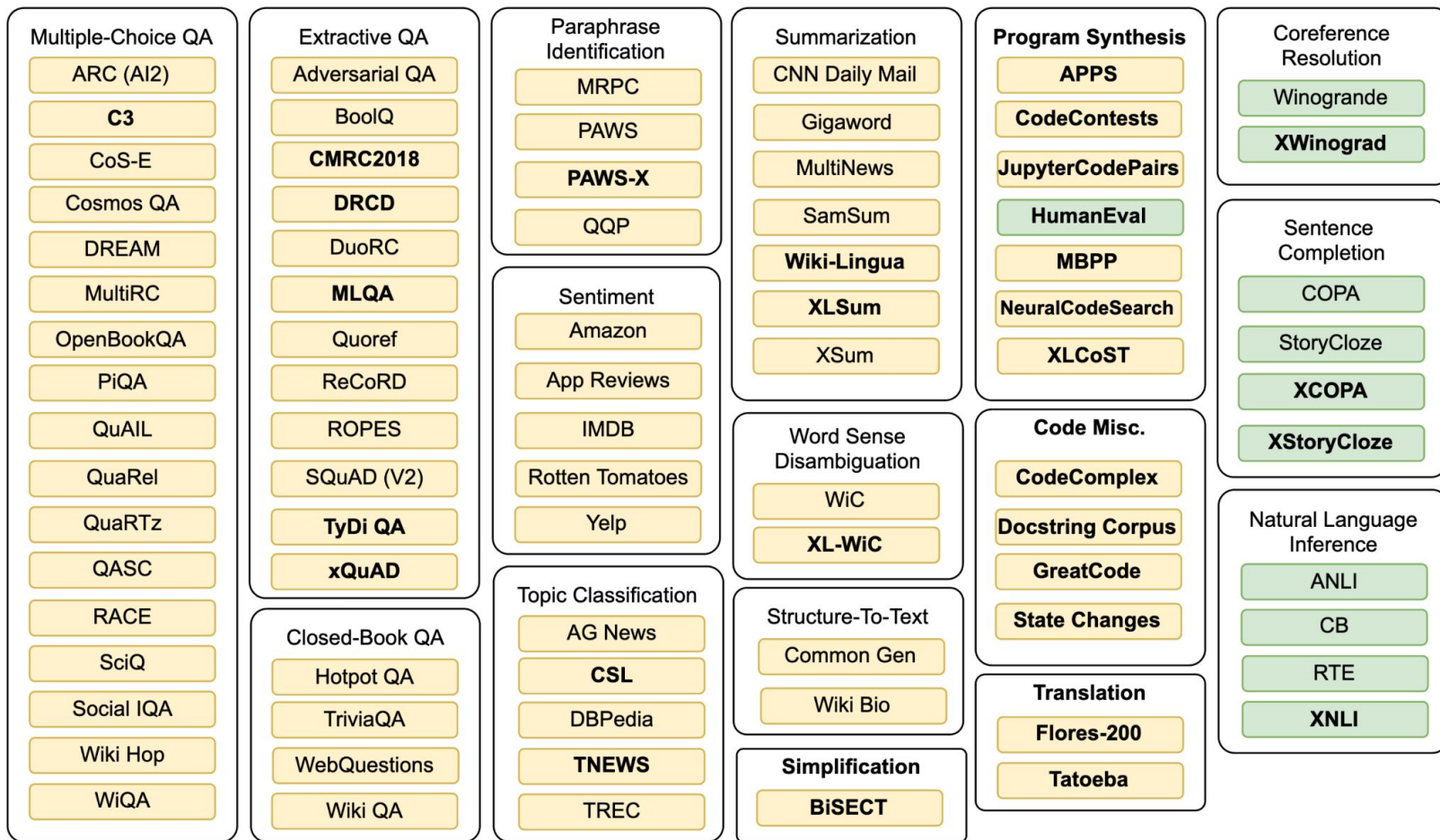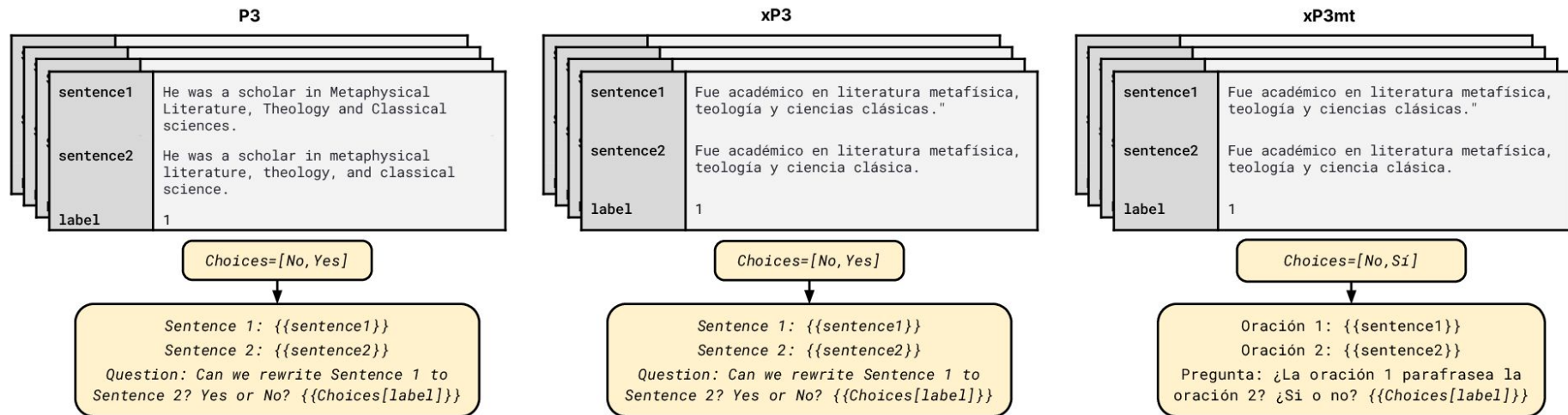*From "Multitask Prompted Training Enables Zero-Shot Task Generalization" by Sanh et al.*

From "Multitask Prompted Training Enables Zero-Shot Task Generalization" by Sanh et al.

## Multiple-Choice QA
- ARC (AI2)
- **C3**
- CoS-E
- Cosmos QA
- DREAM
- MultiRC
- OpenBookQA
- PiQA
- QuAIL
- QuaRel
- QuaRTz
- QASC
- RACE
- SciQ
- Social IQA
- Wiki Hop
- WiQA

## Extractive QA
- Adversarial QA
- BoolQ
- **CMRC2018**
- **DRCD**
- DuoRC
- **MLQA**
- Quoref
- ReCoRD
- ROPES
- SQuAD (V2)
- **TyDi QA**
- **xQuAD**

## Closed-Book QA
- Hotpot QA
- TriviaQA
- WebQuestions
- Wiki QA

## Paraphrase Identification
- MRPC
- PAWS
- **PAWS-X**
- QQP

## Sentiment
- Amazon
- App Reviews
- IMDB
- Rotten Tomatoes
- Yelp

## Topic Classification
- AG News
- **CSL**
- DBPedia
- **TNEWS**
- TREC

## Summarization
- CNN Daily Mail
- Gigaword
- MultiNews
- SamSum
- **Wiki-Lingua**
- **XLSum**
- XSum

## Word Sense Disambiguation
- WiC
- **XL-WiC**

## Structure-To-Text
- Common Gen
- Wiki Bio

## Simplification
- **BiSECT**

## Program Synthesis
- **APPS**
- **CodeContests**
- **JupyterCodePairs**
- HumanEval
- **MBPP**
- NeuralCodeSearch
- **XLCoST**

## Code Misc.
- **CodeComplex**
- **Docstring Corpus**
- **GreatCode**
- **State Changes**

## Translation
- **Flores-200**
- **Tatoeba**

## Coreference Resolution
- Winogrande
- **XWinograd**

## Sentence Completion
- COPA
- StoryCloze
- **XCOPA**
- **XStoryCloze**

## Natural Language Inference
- ANLI
- CB
- RTE
- **XNLI**

*From "Crosslingual Generalization through Multitask Finetuning" by Muennighoff et al.*

**P3**

| sentence1 | He was a scholar in Metaphysical Literature, Theology and Classical sciences. |
| sentence2 | He was a scholar in metaphysical literature, theology, and classical science. |
| label | 1 |

*Choices=[No,Yes]*

*Sentence 1: {{sentence1}}*
*Sentence 2: {{sentence2}}*
*Question: Can we rewrite Sentence 1 to Sentence 2? Yes or No? {{Choices[label]}}*

**xP3**

| sentence1 | Fue académico en literatura metafísica, teología y ciencias clásicas." |
| sentence2 | Fue académico en literatura metafísica, teología y ciencia clásica. |
| label | 1 |

*Choices=[No,Yes]*

*Sentence 1: {{sentence1}}*
*Sentence 2: {{sentence2}}*
*Question: Can we rewrite Sentence 1 to Sentence 2? Yes or No? {{Choices[label]}}*

**xP3mt**

| sentence1 | Fue académico en literatura metafísica, teología y ciencias clásicas." |
| sentence2 | Fue académico en literatura metafísica, teología y ciencia clásica. |
| label | 1 |

*Choices=[No,Sí]*

*Oración 1: {{sentence1}}*
*Oración 2: {{sentence2}}*
*Pregunta: ¿La oración 1 parafrasea la oración 2? ¿Si o no? {{Choices[label]}}*

*From "Crosslingual Generalization through Multitask Finetuning" by Muennighoff et al.*

# Multilingual Multitask Generalization



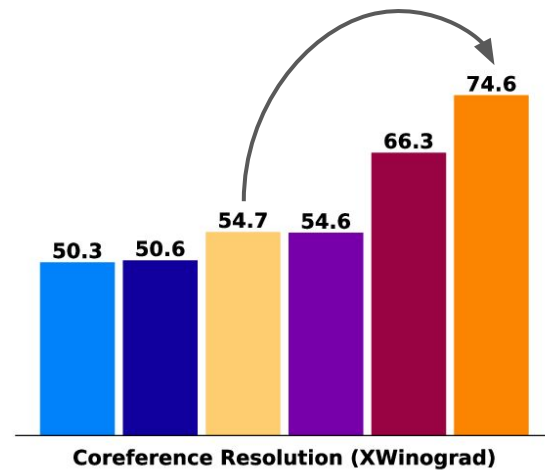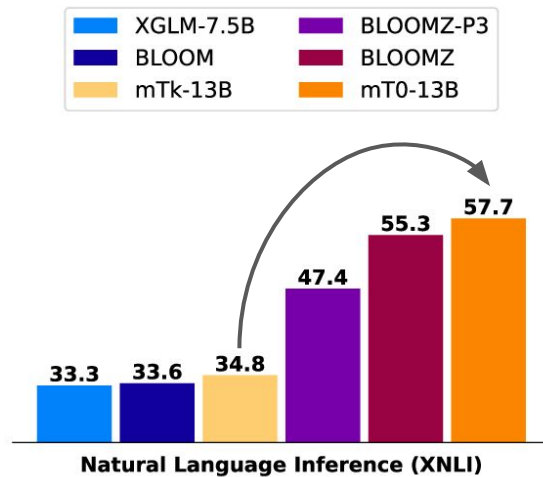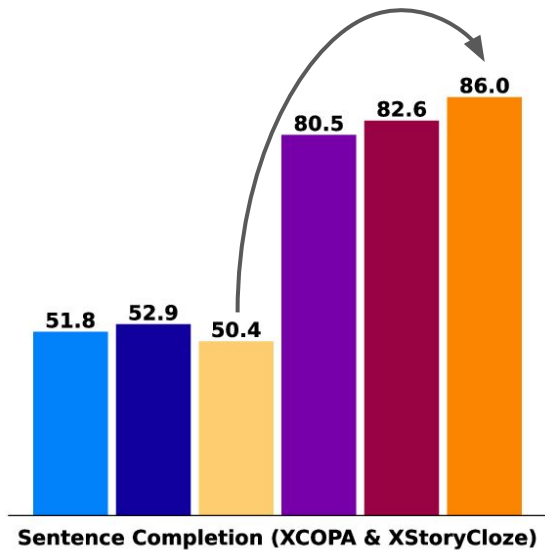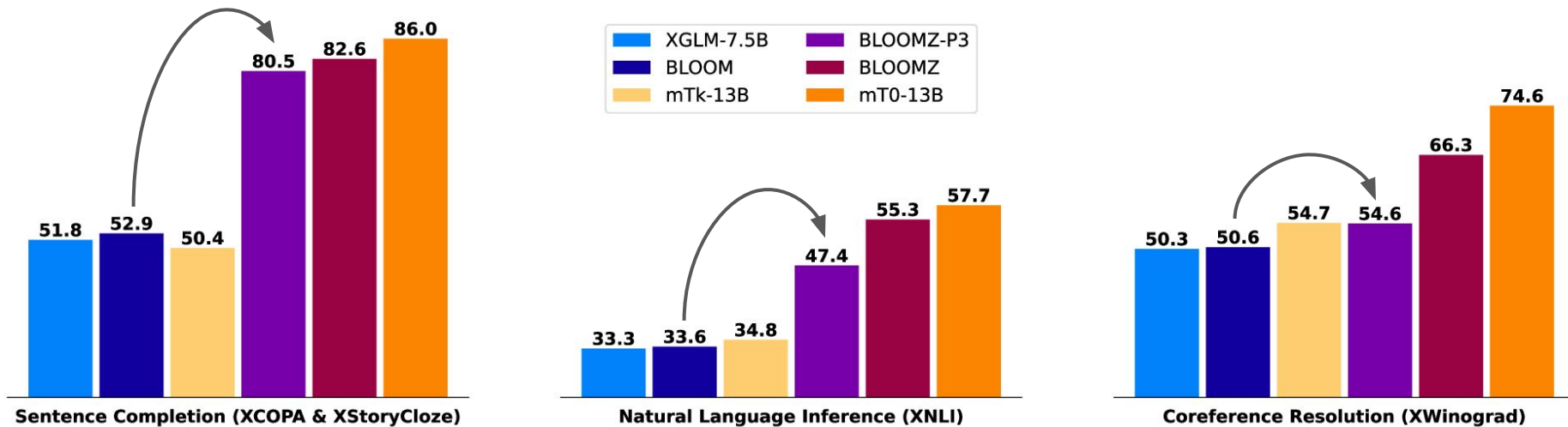**Sentence Completion (XCOPA & XStoryCloze)**
- XGLM-7.5B: 51.8
- BLOOM: 52.9
- mTk-13B: 50.4
- BLOOMZ-P3: 80.5
- BLOOMZ: 82.6
- mT0-13B: 86.0

**Natural Language Inference (XNLI)**
- XGLM-7.5B: 33.3
- BLOOM: 33.6
- mTk-13B: 34.8
- BLOOMZ-P3: 47.4
- BLOOMZ: 55.3
- mT0-13B: 57.7

**Coreference Resolution (XWinograd)**
- XGLM-7.5B: 50.3
- BLOOM: 50.6
- mTk-13B: 54.7
- BLOOMZ-P3: 54.6
- BLOOMZ: 66.3
- mT0-13B: 74.6

*From "Crosslingual Generalization through Multitask Finetuning" by Muennighoff et al.*

Multilingual Multitask Generalization

From "Crosslingual Generalization through Multitask Finetuning" by Muennighoff et al.
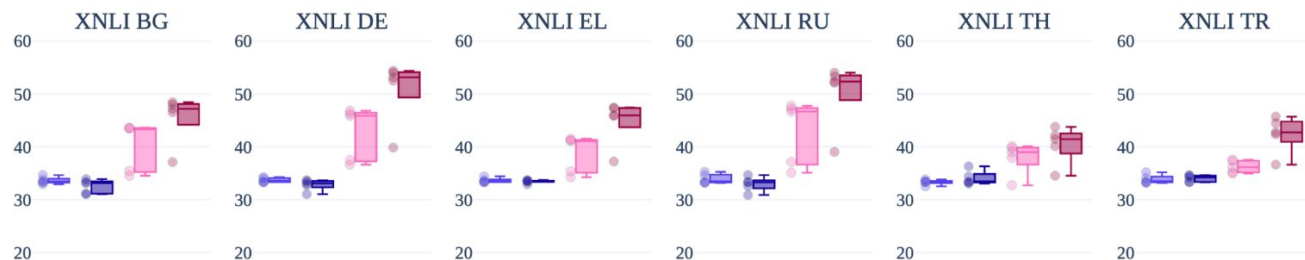
**Multilingual Multitask Generalization**

Legend:
- XGLM-7.5B
- BLOOM
- mTk-13B
- BLOOMZ-P3
- BLOOMZ
- mT0-13B

**Sentence Completion (XCOPA & XStoryCloze)**
- 51.8
- 52.9
- 50.4
- 80.5
- 82.6
- 86.0

**Natural Language Inference (XNLI)**
- 33.3
- 33.6
- 34.8
- 47.4
- 55.3
- 57.7

**Coreference Resolution (XWinograd)**
- 50.3
- 50.6
- 54.7
- 54.6
- 66.3
- 74.6

*From "Crosslingual Generalization through Multitask Finetuning" by Muennighoff et al.*
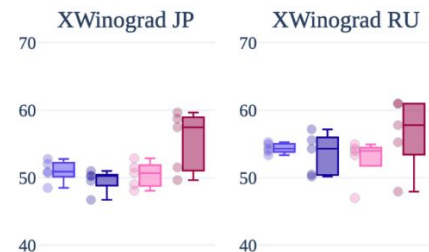
# Multilingual Multitask Generalization

**Sentence Completion (XCOPA & XStoryCloze)**
- XGLM-7.5B: 51.8
- BLOOM: 52.9
- mTk-13B: 50.4
- BLOOMZ-P3: 80.5
- BLOOMZ: 82.6
- mT0-13B: 86.0

**Natural Language Inference (XNLI)**
- XGLM-7.5B: 33.3
- BLOOM: 33.6
- mTk-13B: 34.8
- BLOOMZ-P3: 47.4
- BLOOMZ: 55.3
- mT0-13B: 57.7

**Coreference Resolution (XWinograd)**
- XGLM-7.5B: 50.3
- BLOOM: 50.6
- mTk-13B: 54.7
- BLOOMZ-P3: 54.6
- BLOOMZ: 66.3
- mT0-13B: 74.6

Legend: XGLM-7.5B, BLOOM, mTk-13B, BLOOMZ-P3, BLOOMZ, mT0-13B

*From "Crosslingual Generalization through Multitask Finetuning" by Muennighoff et al.*

**Multilingual Multitask Generalization**

Legend:
- XGLM-7.5B
- BLOOM
- mTk-13B
- BLOOMZ-P3
- BLOOMZ
- mT0-13B

**Sentence Completion (XCOPA & XStoryCloze)**
- 51.8
- 52.9
- 50.4
- 80.5
- 82.6
- 86.0

**Natural Language Inference (XNLI)**
- 33.3
- 33.6
- 34.8
- 47.4
- 55.3
- 57.7

**Coreference Resolution (XWinograd)**
- 50.3
- 50.6
- 54.7
- 54.6
- 66.3
- 74.6

*From "Crosslingual Generalization through Multitask Finetuning" by Muennighoff et al.*

# Performance on languages that were never intentionally trained on

**Natural Language Inference**

**Coreference Resolution**

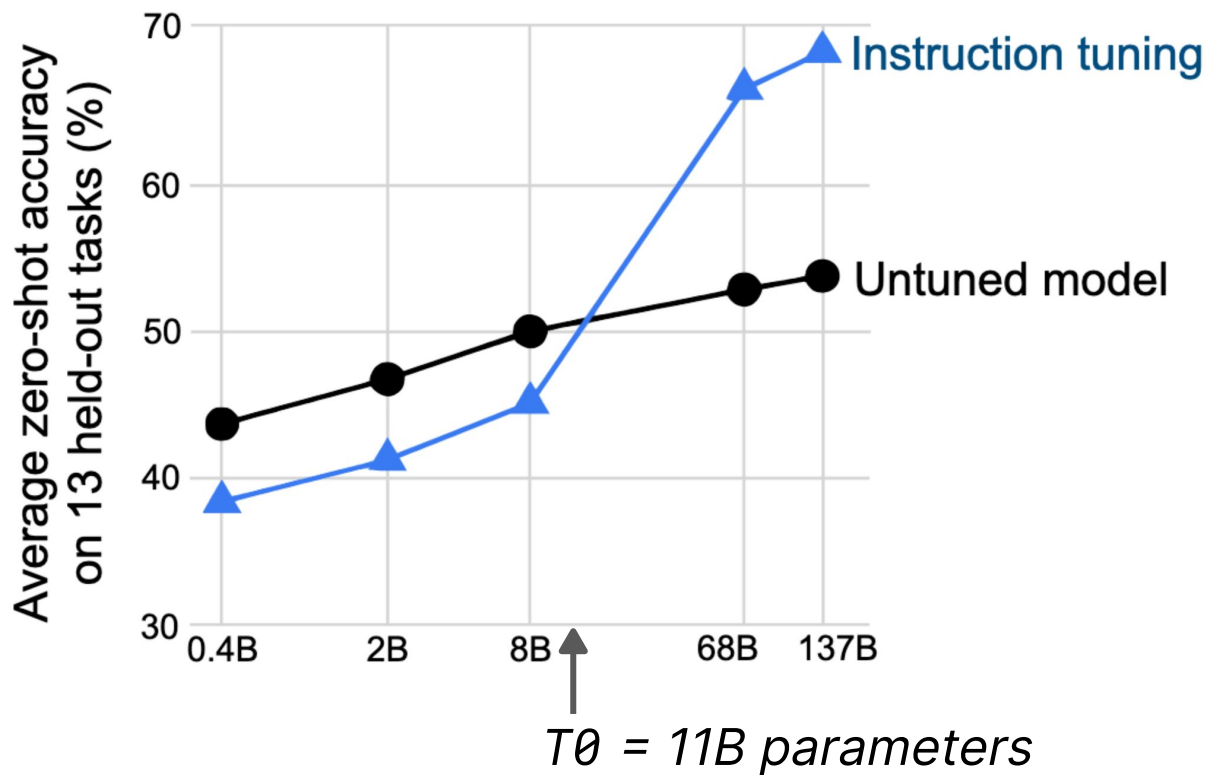XNLI BG · XNLI DE · XNLI EL · XNLI RU · XNLI TH · XNLI TR

XWinograd JP · XWinograd RU

**Sentence Completion**

XCOPA ET · XCOPA HT · XCOPA IT · XCOPA QU · XCOPA TH · XCOPA TR · XStoryCloze MY · XStoryCloze RU

BLOOM-7.1B   BLOOM   BLOOMZ-7.1B   BLOOMZ

*From "Crosslingual Generalization through Multitask Finetuning" by Muennighoff et al.*
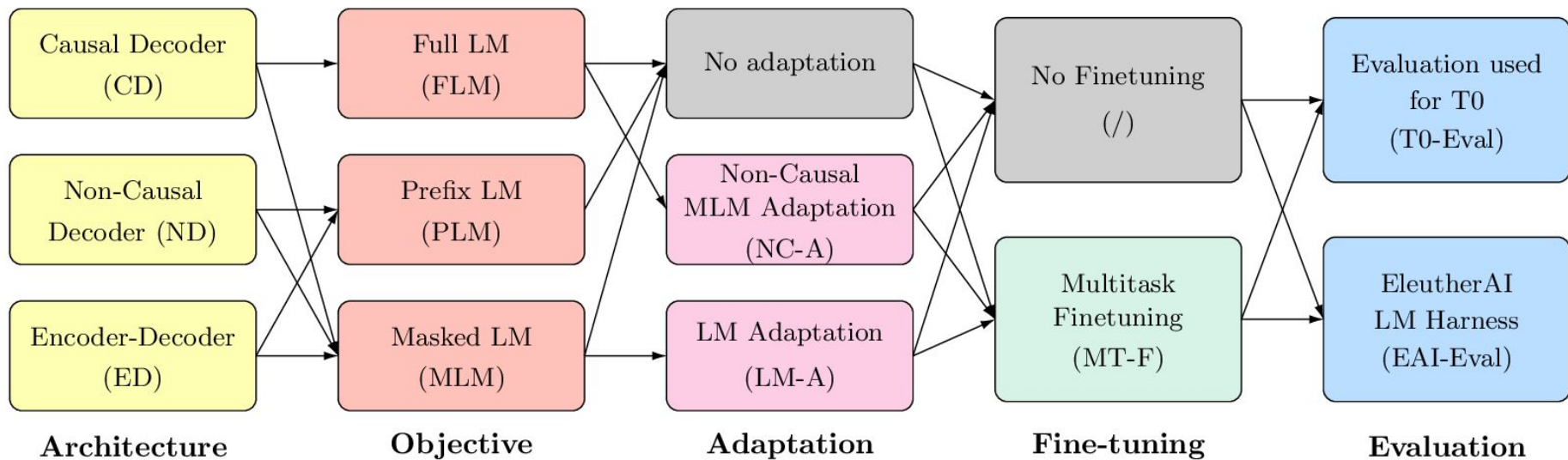
From "Crosslingual Generalization through Multitask Finetuning" by Muennighoff et al.

Performance on **_<u>held-out</u>_** tasks

*From "Fine-Tuned Language Models are Zero-Shot Learners" by Wei et al.*

| Architecture | Objective | Adaptation | Fine-tuning | Evaluation |
|---|---|---|---|---|
| Causal Decoder (CD) | Full LM (FLM) | No adaptation | No Finetuning (/) | Evaluation used for T0 (T0-Eval) |
| Non-Causal Decoder (ND) | Prefix LM (PLM) | Non-Causal MLM Adaptation (NC-A) | Multitask Finetuning (MT-F) | EleutherAI LM Harness (EAI-Eval) |
| Encoder-Decoder (ED) | Masked LM (MLM) | LM Adaptation (LM-A) | | |

*From "What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization?" by Wang et al.*

Causal Decoder · Non-causal Decoder · Encoder-Decoder
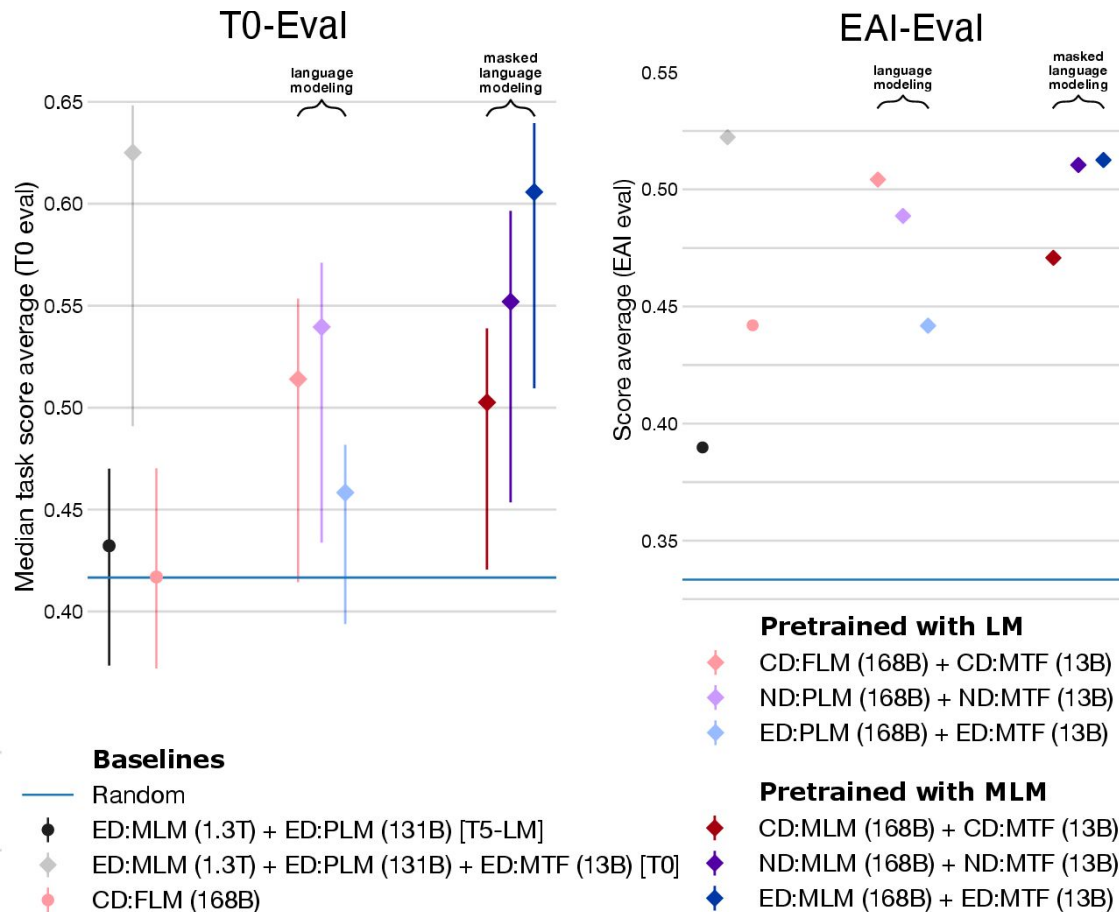
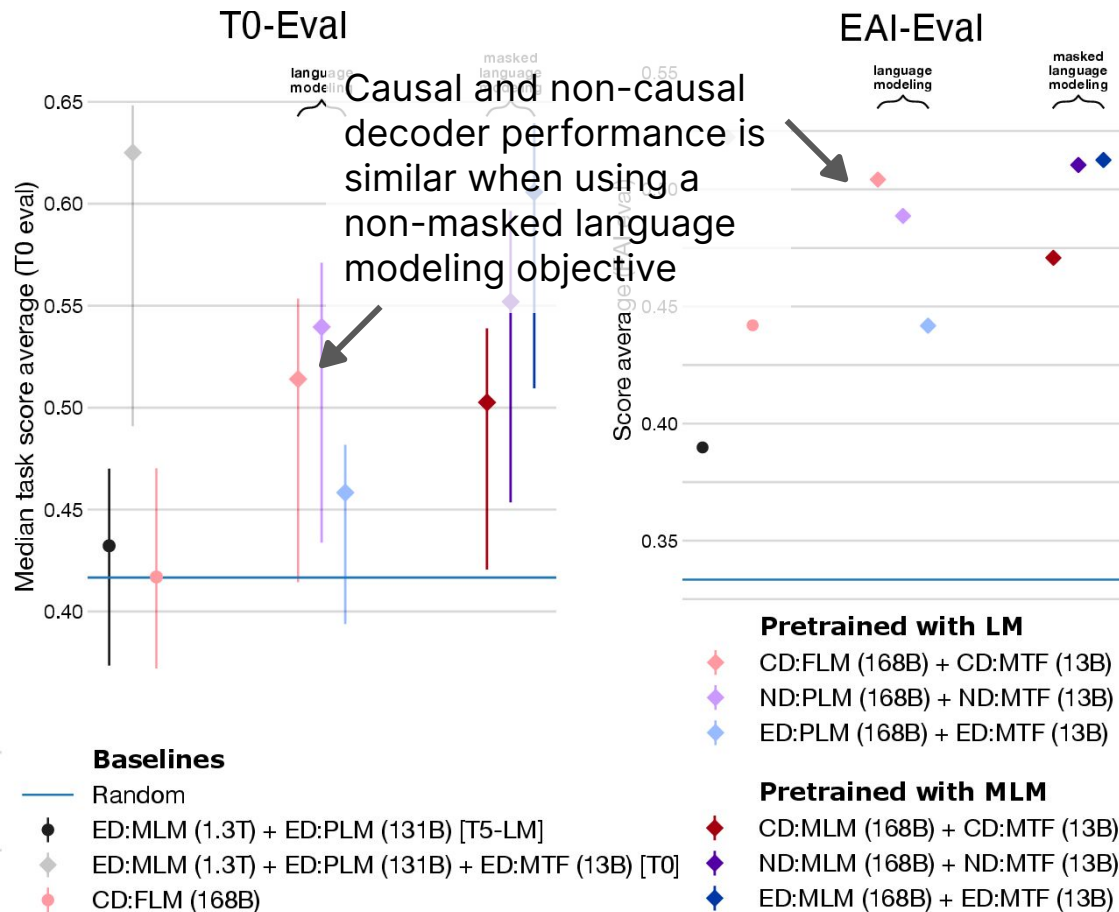**Full Language Modeling**   May [targets] the force be with you

**Prefix Language Modeling**   May the force [targets] be with you

**Masked Language Modeling**   May [targets] the force be with you

*From "What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization?" by Wang et al.*

**T0-Eval**

Median task score average (T0 eval)

language modeling

masked language modeling

0.65
0.60
0.55
0.50
0.45
0.40

**Baselines**
— Random
● ED:MLM (1.3T) + ED:PLM (131B) [T5-LM]
◆ ED:MLM (1.3T) + ED:PLM (131B) + ED:MTF (13B) [T0]
● CD:FLM (168B)

**EAI-Eval**

Score average (EAI eval)

language modeling

masked language modeling

0.55
0.50
0.45
0.40
0.35

**Pretrained with LM**
◆ CD:FLM (168B) + CD:MTF (13B)
◆ ND:PLM (168B) + ND:MTF (13B)
◆ ED:PLM (168B) + ED:MTF (13B)

**Pretrained with MLM**
◆ CD:MLM (168B) + CD:MTF (13B)
◆ ND:MLM (168B) + ND:MTF (13B)
◆ ED:MLM (168B) + ED:MTF (13B)

*From "What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization?" by Wang et al.*

T0-Eval

EAI-Eval

Causal and non-causal decoder performance is similar when using a non-masked language modeling objective

**Baselines**
— Random
● ED:MLM (1.3T) + ED:PLM (131B) [T5-LM]
◆ ED:MLM (1.3T) + ED:PLM (131B) + ED:MTF (13B) [T0]
● CD:FLM (168B)

**Pretrained with LM**
◆ CD:FLM (168B) + CD:MTF (13B)
◆ ND:PLM (168B) + ND:MTF (13B)
◆ ED:PLM (168B) + ED:MTF (13B)

**Pretrained with MLM**
◆ CD:MLM (168B) + CD:MTF (13B)
◆ ND:MLM (168B) + ND:MTF (13B)
◆ ED:MLM (168B) + ED:MTF (13B)

*From "What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization?" by Wang et al.*

## T0-Eval

## EAI-Eval

Encoder-decoder does very poorly when using a non-masked language modeling objective

**Baselines**
— Random
● ED:MLM (1.3T) + ED:PLM (131B) [T5-LM]
◆ ED:MLM (1.3T) + ED:PLM (131B) + ED:MTF (13B) [T0]
● CD:FLM (168B)

**Pretrained with LM**
◆ CD:FLM (168B) + CD:MTF (13B)
◆ ND:PLM (168B) + ND:MTF (13B)
◆ ED:PLM (168B) + ED:MTF (13B)

**Pretrained with MLM**
◆ CD:MLM (168B) + CD:MTF (13B)
◆ ND:MLM (168B) + ND:MTF (13B)
◆ ED:MLM (168B) + ED:MTF (13B)

*From "What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization?" by Wang et al.*

For masked language modeling, the opposite is true - non-causal visibility on the prefix helps a lot!

*From "What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization?" by Wang et al.*

**Zero-shot**

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```
1   Translate English to French:        ←  task description

2   cheese =>  .....................    ←  prompt
```
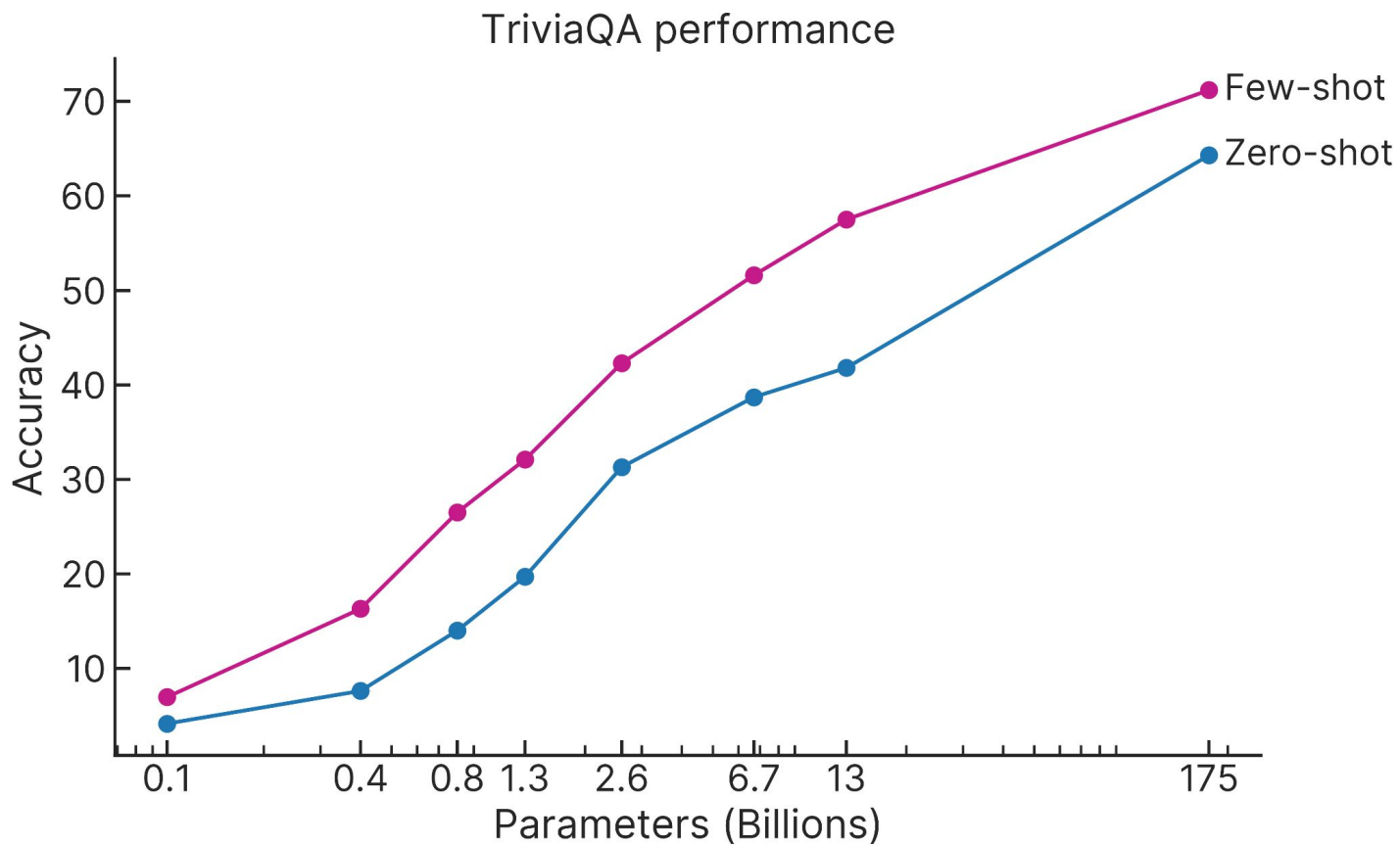
**Few-shot**

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1   Translate English to French:        ←  task description

2   sea otter => loutre de mer          ←  examples

3   peppermint => menthe poivrée

4   plush girafe => girafe peluche

5   cheese =>  .....................    ←  prompt
```
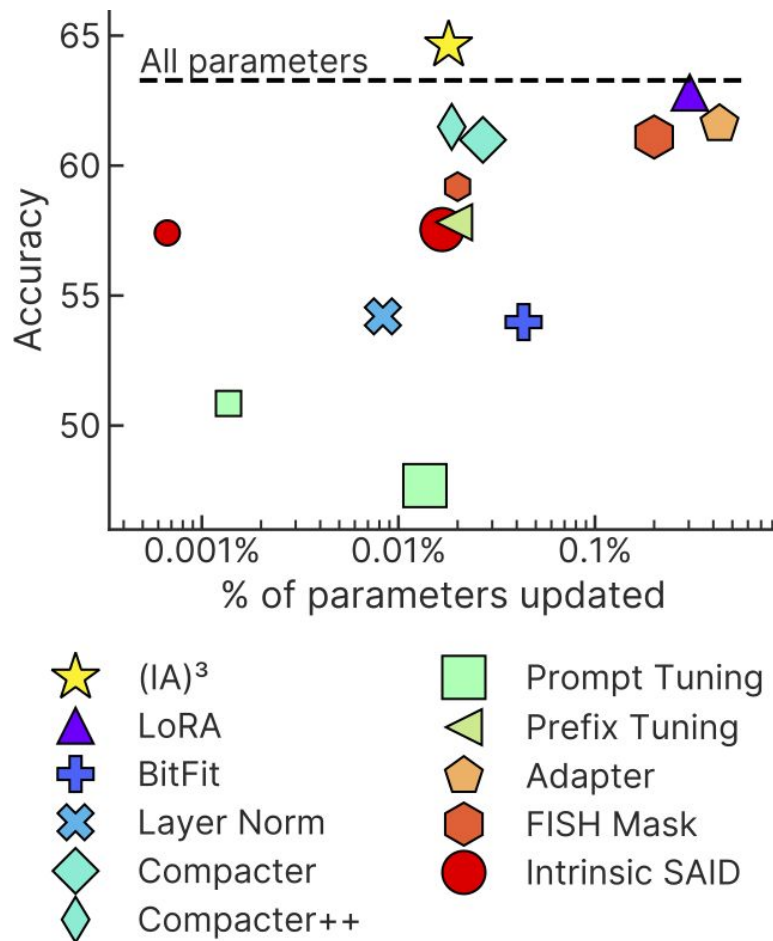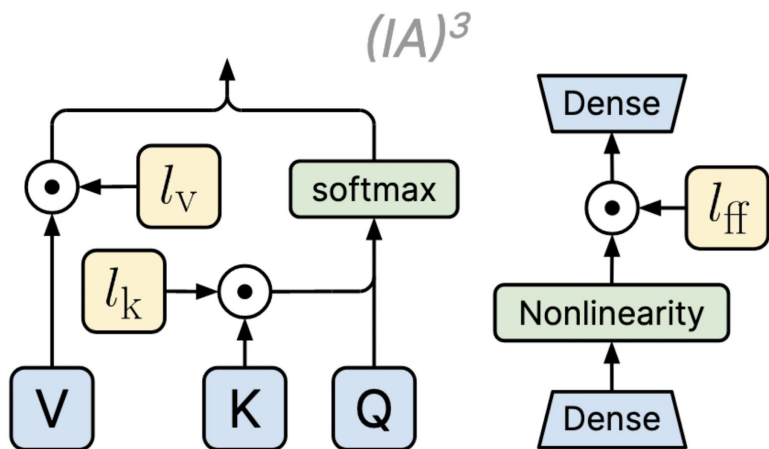
**Fine-tuning**

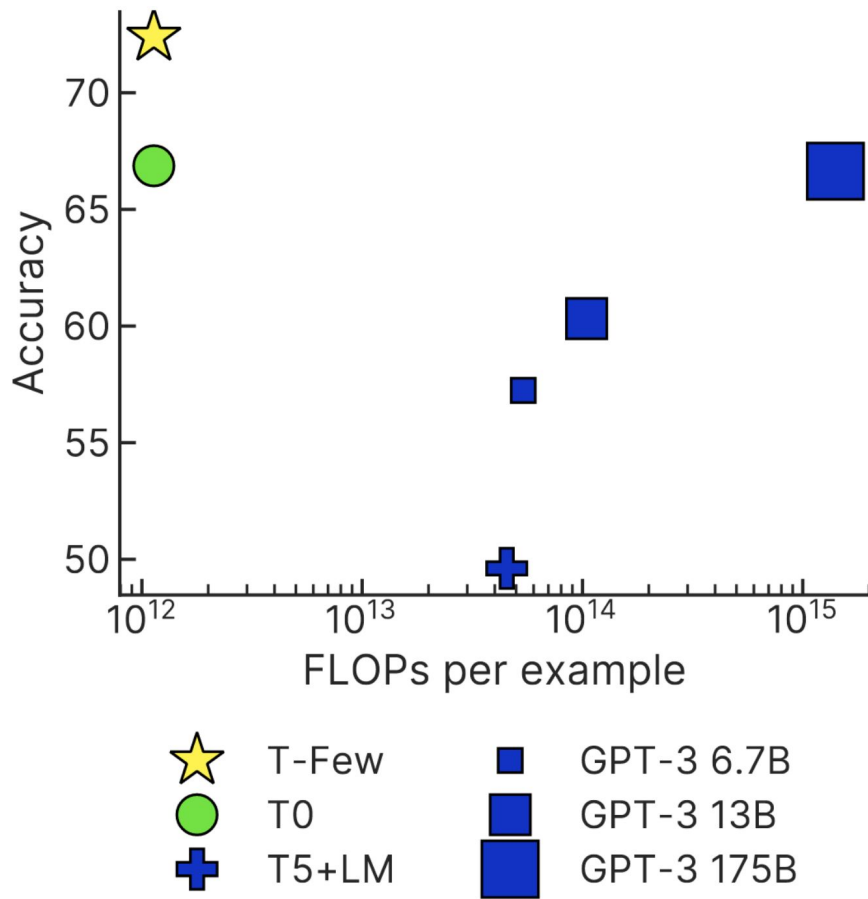The model is trained via repeated gradient updates using a large corpus of example tasks.

```
1   sea otter => loutre de mer          ←  example #1
```
↓
**gradient update**
↓
```
1   peppermint => menthe poivrée        ←  example #2
```
↓
**gradient update**
↓
● ● ●
↓
```
1   plush giraffe => girafe peluche     ←  example #N
```

**gradient update**

```
1   cheese =>  .....................    ←  prompt
```

*From "Language Models are Few-Shot Learners" by Brown et al.*

TriviaQA performance

From "Language Models are Few-Shot Learners" by Brown et al.

$(IA)^3$

From "Few-Shot Parameter-Efficient Fine-Tuning is Better and Cheaper than In-Context Learning", Liu et al. 2022

| Method | Inference FLOPs | Training FLOPs | Disk space |
|---|---|---|---|
| T-Few | 1.1e12 | 2.7e16 | 4.2 MB |
| T0 [1] | 1.1e12 | 0 | 0 B |
| T5+LM [14] | 4.5e13 | 0 | 16 kB |
| GPT-3 6.7B [4] | 5.4e13 | 0 | 16 kB |
| GPT-3 13B [4] | 1.0e14 | 0 | 16 kB |
| GPT-3 175B [4] | 1.4e15 | 0 | 16 kB |

*From "Few-Shot Parameter-Efficient Fine-Tuning is Better and Cheaper than In-Context Learning", Liu et al. 2022*

| Method | Acc. |
|---|---|
| T-Few | 75.8% |
| Human baseline [2] | 73.5% |
| PET [50] | 69.6% |
| SetFit [51] | 66.9% |
| GPT-3 [4] | 62.7% |

Table 2: Top-5 best methods on RAFT as of writing. T-Few is the first method to outperform the human baseline and achieves over 6% higher accuracy than the next-best method.

Legend:
- ★ T-Few
- ● T0
- ✚ T5+LM
- ■ GPT-3 6.7B
- ■ GPT-3 13B
- ■ GPT-3 175B

*From "Few-Shot Parameter-Efficient Fine-Tuning is Better and Cheaper than In-Context Learning", Liu et al. 2022*

# References

[Multitask Prompted Training Enables Zero-Shot Task Generalization](#)

[Crosslingual Generalization through Multitask Finetuning](#)

[What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization?](#)

[Few-Shot Parameter-Efficient Fine-Tuning is Better and Cheaper than In-Context Learning](#)

[BLOOM: A 176B-Parameter Open-Access Multilingual Language Model](#)

Please give me feedback:

[http://bit.ly/colin-talk-feedback](http://bit.ly/colin-talk-feedback)

Thanks!