# Multitask Training and Instruction-Based Fine-Tuning

## Flan T5/ T0

**Groupe_5**

Abdelaziz Guelfane
Imane Meziany
Malek Bouhadida
Mohammed El Barhichi
Yousra Yakhou

CentraleSupélec

- Pre-2018 : **Task-Specific Training**
  - RNNs, LSTMs, GRUs were trained on specific tasks.
  - Even Seq2Seq models with Attention (Transformer architecture) involved training a separate model for each task.

- Pre-2018 : **Task-Specific Training**
  - RNNs, LSTMs, GRUs were trained on specific tasks.
  - Even Seq2Seq models with Attention (Transformer architecture) involved training a separate model for each task.

- 2019 : **Text-To-Text framework**
  - T5 reformulates any NLP task to a text generation problem.
  - Example : Instead of predicting a class label, the model generates the label as text (e.g., "positive sentiment").
  - However, fine-tuning on specific tasks was still needed for optimal performance.

# Timeline - Paradigm Change

- Pre-2018 : **Task-Specific Training**
  - RNNs, LSTMs, GRUs were trained on specific tasks
  - Even Seq2Seq models with Attention (Transformer architecture) involved training a separate model for each task.

- 2019 : **Text-To-Text framework**
  - T5 reformulates any NLP task to a text generation problem.
  - Example : Instead of predicting a class label, the model generates the label as text (e.g., "positive sentiment").
  - However, fine-tuning on specific tasks was still needed for optimal performance.

- 2020-2022 : **Instruction Tuning on Multi-tasks**
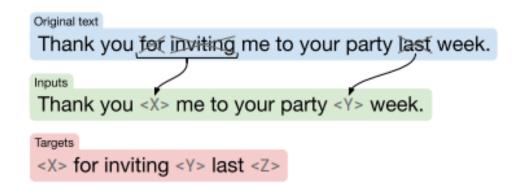  - **T0** : Sanh, V., A. Webson, C. Raffel, et al. "**Multitask Prompted Training Enables Zero-Shot Task Generalization.**" In ICLR, 2022.
  - **Flan-T5** : Chung, H.W., L. Hou, S. Longpre, et al. "**Scaling Instruction-Finetuned Language Models.**" In CoRR, 2022.

- Pre-2018 : **Task-Specific Training**
  - RNNs, LSTMs, GRUs were trained on specific tasks
  - Even Seq2Seq models with Attention (Transformer architecture) involved training a separate model for each task.

- 2019 : **Text-To-Text framework**
  - T5 reformulates any NLP task to a text generation problem.
  - Example : Instead of predicting a class label, the model generates the label as text (e.g., "positive sentiment").
  - However, fine-tuning on specific tasks was still needed for optimal performance.

  **Better Generalization ?**

- 2020-2022 : **Instruction Tuning on Multi-tasks**
  - **T0** : Sanh, V., A. Webson, C. Raffel, et al. "**Multitask Prompted Training Enables Zero-Shot Task Generalization.**" In ICLR, 2022.
  - **Flan-T5** : Chung, H.W., L. Hou, S. Longpre, et al. "**Scaling Instruction-Finetuned Language Models.**" In CoRR, 2022.

# T5: Text-to-Text Transfer Transformer

- **T5** reformulates any NLP task as a sequence-to-sequence generation problem.
- It is based on an **Encoder-Decoder** architecture and trained using **span corruption.**
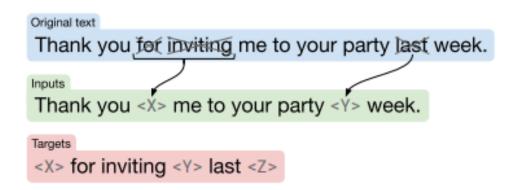- It was pretained on **C4 (Colossal Clean Crawled Corpus).**

Original text
Thank you for inviting me to your party last week.

Inputs
Thank you <X> me to your party <Y> week.

Targets
<X> for inviting <Y> last <Z>

- **T5+LM** is the T5 version designed for language modeling tasks.
- It is T5 trained on 100B additional tokens from C4 on a **causal language modeling** objective.

The model predicts the next token in a sequence, given all prior tokens.

# T5: Text-to-Text Transfer Transformer

- **T5** reformulates any NLP task as a sequence-to-sequence generation problem.
- It is based on an **Encoder-Decoder** architecture and trained using **span corruption.**
- It was pretained on **C4 (Colossal Clean Crawled Corpus).**

Original text
Thank you for inviting me to your party last week.

Inputs
Thank you <X> me to your party <Y> week.

Targets
<X> for inviting <Y> last <Z>

- **T5+LM** is the T5 version designed for language modeling tasks.
- It is T5 trained on 100B additional tokens from C4 on a **causal language modeling** objective.

The model predicts the next token in a sequence, given all prior tokens.

- T5-LM, and LLMs in general, attain reasonable zero-shot generalization on a diverse set of tasks.

## Hypothesis :

>>>>> **Consequence of <u>implicit multitask learning</u> in language models' pretraining.**

**Implicit multitask learning**

- Many websites contain lists of trivia Q&As
- >>>>> **Supervised training data** for the task of closedbook question answering.

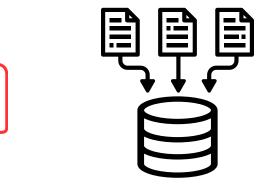1. What does "www" stand for in a website browser?

**Answer:** World Wide Web

2. How long is an Olympic swimming pool (in meters)?

**Answer:** 50 meters

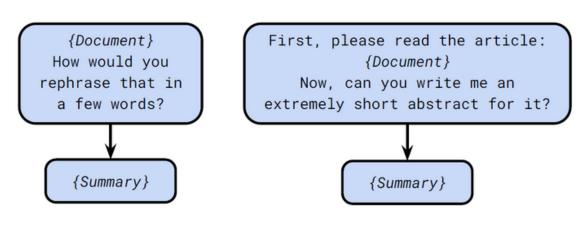3. What countries made up the original Axis powers in World War II?

**Answer:** Germany, Italy, and Japan

# T0: MultiTask Prompted Training Enables Zero-Shot Task Generalization

## Implicit multitask learning

## Limitations..

- Many websites contain lists of trivia Q&As
- >>>>> **Supervised training data** for the task of closedbook question answering.

1. What does "www" stand for in a website browser?
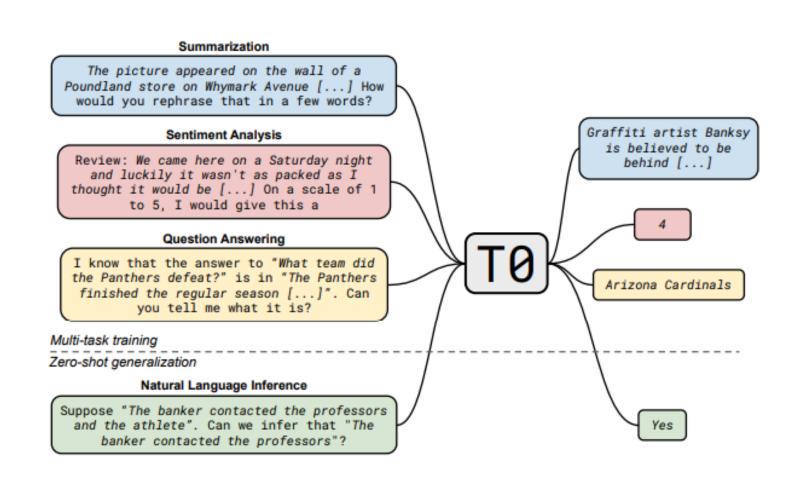
**Answer:** World Wide Web

2. How long is an Olympic swimming pool (in meters)?

**Answer:** 50 meters

3. What countries made up the original Axis powers in World War II?
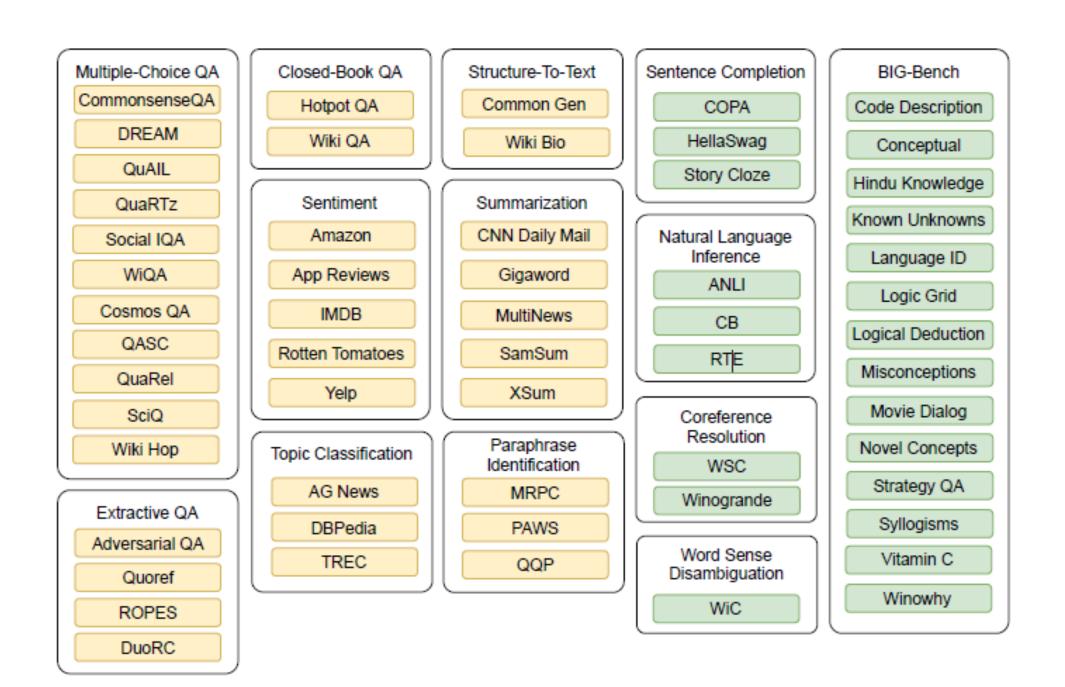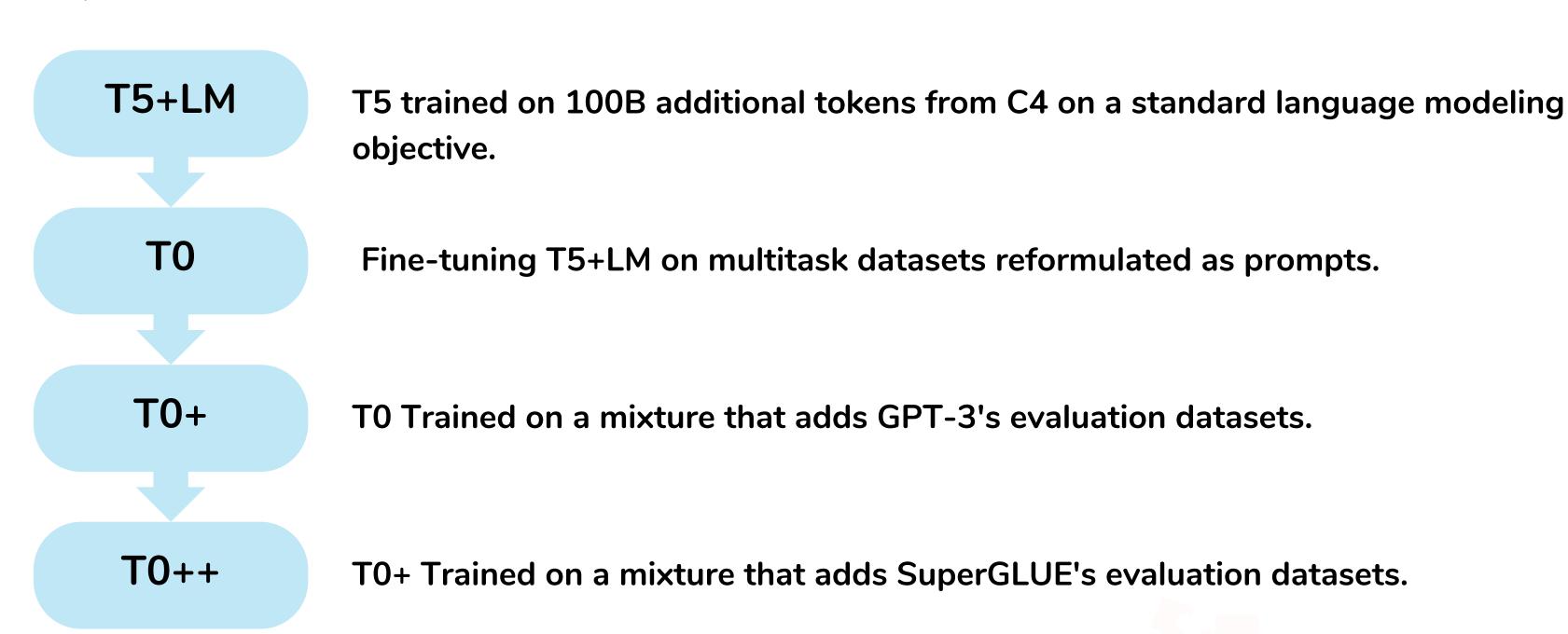
**Answer:** Germany, Italy, and Japan

**1**

**Requires a sufficiently large model + larger corpus of data**

**2**

```
{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}
```

**Sensitive to the wording of prompts**

# T0: MultiTask Prompted Training Enables Zero-Shot Task Generalization

## Implicit multitask learning

- Many websites contain lists of trivia Q&As
- **>>>>** **Supervised training data** for the task of closedbook question answering.

1. What does "www" stand for in a website browser?

Answer: World Wide Web

2. How long is an Olympic swimming pool (in meters)?

Answer: 50 meters

3. What countries made up the original Axis powers in World War II?

Answer: Germany, Italy, and Japan

**VS**

## Explicit multitask learning



**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"?

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

T0

4

Arizona Cardinals

Yes

CentraleSupélec

- **T0** datasets and task taxonomy (T0+ and T0++ are trained on additional datasets).

  - **12 tasks and 62 datasets** with publicly contributed prompts.
  - **Yellow datasets** are in the training mixture.
  - **Green datasets** are held out and represent tasks that were not seen during training.
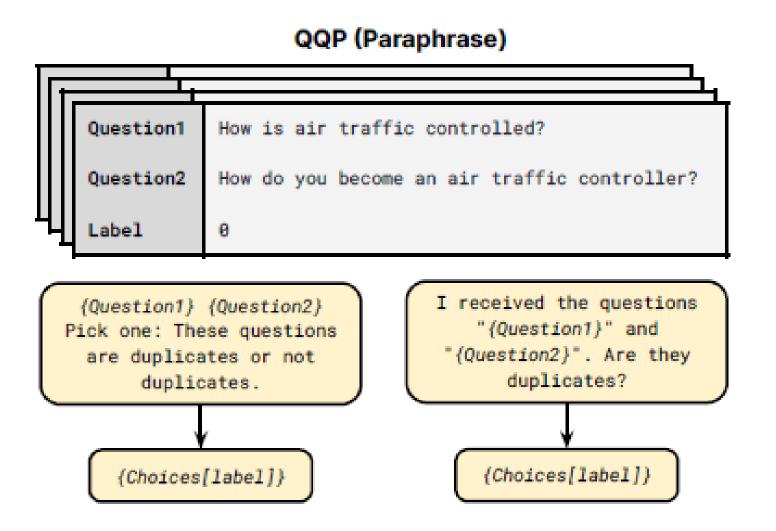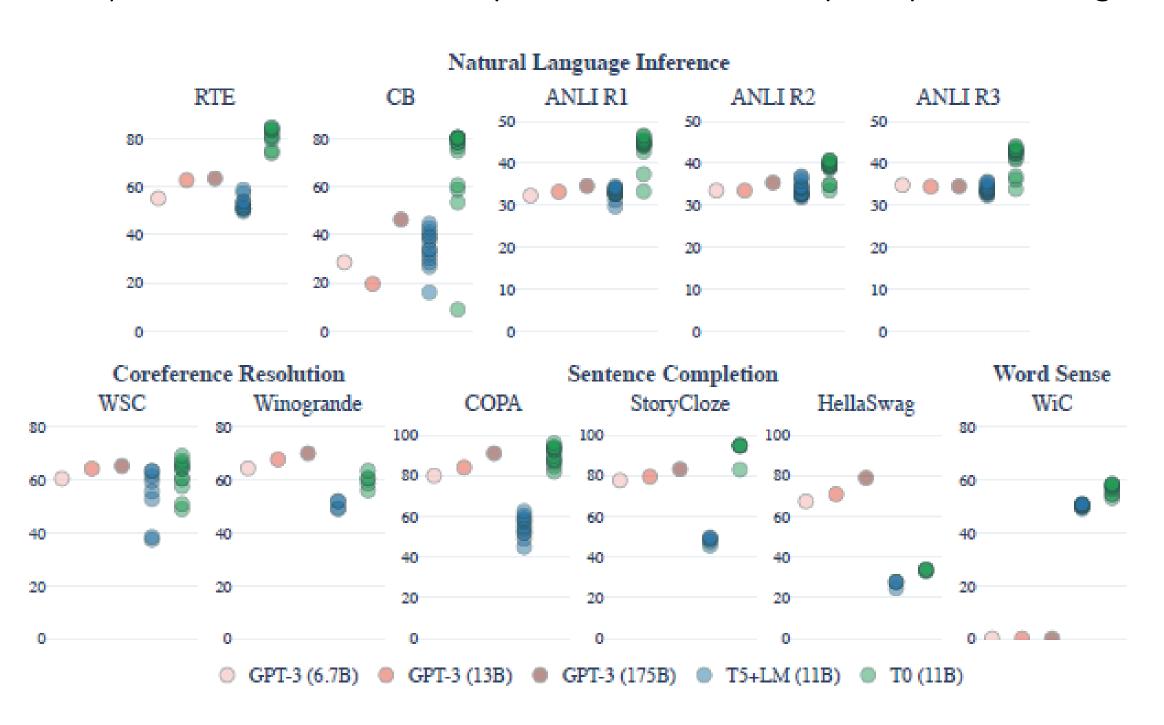


**Multiple-Choice QA**
- CommonsenseQA
- DREAM
- QuAIL
- QuaRTz
- Social IQA
- WiQA
- Cosmos QA
- QASC
- QuaRel
- SciQ
- Wiki Hop

**Extractive QA**
- Adversarial QA
- Quoref
- ROPES
- DuoRC

**Closed-Book QA**
- Hotpot QA
- Wiki QA

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

**Topic Classification**
- AG News
- DBPedia
- TREC

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

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

**Paraphrase Identification**
- MRPC
- PAWS
- QQP

**Sentence Completion**
- COPA
- HellaSwag
- Story Cloze

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

**Coreference Resolution**
- WSC
- Winogrande

**Word Sense Disambiguation**
- WiC

**BIG-Bench**
- Code Description
- Conceptual
- Hindu Knowledge
- Known Unknowns
- Language ID
- Logic Grid
- Logical Deduction
- Misconceptions
- Movie Dialog
- Novel Concepts
- Strategy QA
- Syllogisms
- Vitamin C
- Winowhy

- A quick overview of the models :

**T5+LM** — T5 trained on 100B additional tokens from C4 on a standard language modeling objective.

↓

**T0** — Fine-tuning T5+LM on multitask datasets reformulated as prompts.

↓

**T0+** — T0 Trained on a mixture that adds GPT-3's evaluation datasets.

↓

**T0++** — T0+ Trained on a mixture that adds SuperGLUE's evaluation datasets.

≫≫≫ The above **T0 variants** are all initialized from the **11B parameters** version of **T5+LM**.

CentraleSupélec

**Models learn to understand the prompts as task instructions which help them generalize to held-out tasks**



**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.

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

{Choices[label]}

{Choices[label]}

**How did generalization to held-out tasks improve?**

- Results for **T0** task generalization experiments compared to **GPT-3**.
- T5+LM (baseline model) is the same as T0 except without multitask prompted training.



**Better Generalization !**

Each dot is the performance of one evaluation prompt.

CentraleSupélec

- Results for a subset of BIG-bench models compared to T0, T0+ and T0++.

**Better Generalization !**

- So far, we have seen that..

**》》》》》**    **Models learn to understand the prompts as task instructions which help them generalize to held-out tasks.**

**》》》》》**    **Multitask prompted training improve generalization to held-out tasks**

- Now,

**》》》》》**    **Does training on a wider range of prompts improve robustness to prompt wording?**

- Zero-shot performance of T0 and T5+LM when increasing number of training prompts per dataset.



Each dot is the performance of one evaluation prompt.
Adding more training prompts consistently leads to higher median performance and generally lower interquartile range for held-out tasks.

**What happens if we scale even further?**

- Increasing the **number of tasks** and **data diversity**?
- Leveraging **larger models**?
- Incorporating **reasoning capabilities** like Chain-of-Thought (CoT)?

**What happens if we scale even further?**

- Increasing the **number of tasks** and **data diversity**?
- Leveraging **larger models**?
- Incorporating **reasoning capabilities** like Chain-of-Thought (CoT)?

**1**

**Scaling Tasks and Datasets**

# Flan-T5 : Scaling Instruction-Finetuned Language Models

**1** Scaling Tasks and Datasets

**+**

**2** Adding Chain-of-Thought Reasoning

# Flan-T5 : Scaling Instruction-Finetuned Language Models



**1** Scaling Tasks and Datasets

**+**

**2** Adding Chain-of-Thought Reasoning

**+**

**3** Using Larger Language Models

- Dataset used :

- Scaling Finetuning Tasks :

  **473 datasets**

  **146 task categories**

  **1 836 Tasks !**



**Finetuning tasks**

**TO-SF**

Commonsense reasoning
Question generation
Closed-book QA
Adversarial QA
Extractive QA
Title/context generation
Topic classification
Struct-to-text
...

*55 Datasets, 14 Categories, 193 Tasks*

**Muffin**

Natural language inference      Closed-book QA
Code instruction gen.            Conversational QA
Program synthesis                Code repair
Dialog context generation        ...

*69 Datasets, 27 Categories, 80 Tasks*

**CoT (Reasoning)**

Arithmetic reasoning            Explanation generation
Commonsense Reasoning           Sentence composition
Implicit reasoning              ...

*9 Datasets, 1 Category, 9 Tasks*

**Natural Instructions v2**

Cause effect classification
Commonsense reasoning
Named entity recognition
Toxic language detection
Question answering
Question generation
Program execution
Text categorization
...

*372 Datasets, 108 Categories, 1554 Tasks*

❖ A **Dataset** is an original data source (e.g. SQuAD).
❖ A **Task Category** is unique task setup (e.g. the SQuAD dataset is configurable for multiple task categories such as extractive question answering, query generation, and context generation).
❖ A **Task** is a unique <dataset, task category> pair, with any number of templates which preserve the task category (e.g. query generation on the SQuAD dataset.)

- Tasks & Instructions variety :

Finetuning includes more
diverse tasks :
- zero-shot
- few-shot
- chain-of-thought

- Tasks & Instructions variety :

Finetuning includes more diverse tasks :
- zero-shot
- few-shot
- chain-of-thought

Without chain-of-thought

Instruction without exemplars

Answer the following yes/no question.

Can you write a whole Haiku in a single tweet?

⟹ yes

- Tasks & Instructions variety :

Finetuning includes more diverse tasks :
- zero-shot
- few-shot
- chain-of-thought



Without chain-of-thought

Instruction without exemplars

Answer the following yes/no question.

Can you write a whole Haiku in a single tweet?

⟶ yes

Instruction with exemplars

Q: Answer the following yes/no question.
Could a dandelion suffer from hepatitis?
A: no

Q: Answer the following yes/no question.
Can you write a whole Haiku in a single tweet?
A:

⟶ yes

- Tasks & Instructions variety :

Finetuning includes more diverse tasks :
- zero-shot
- few-shot
- chain-of-thought



Without chain-of-thought

**Instruction without exemplars**

Answer the following yes/no question.

Can you write a whole Haiku in a single tweet?

→ yes

**Instruction with exemplars**

Q: Answer the following yes/no question.
Could a dandelion suffer from hepatitis?
A: no

Q: Answer the following yes/no question.
Can you write a whole Haiku in a single tweet?
A:

→ yes

With chain-of-thought

Answer the following yes/no question by reasoning step-by-step.

Can you write a whole Haiku in a single tweet?

→ A haiku is a japanese three-line poem. That is short enough to fit in 280 characters. The answer is yes.

Q: Answer the following yes/no question by reasoning step-by-step.
Could a dandelion suffer from hepatitis?
A: Hepatitis only affects organisms with livers. Dandelions don't have a liver. The answer is no.

Q: Answer the following yes/no question by reasoning step-by-step.
Can you write a whole Haiku in a single tweet?
A:

→ A haiku is a japanese three-line poem. That is short enough to fit in 280 characters. The answer is yes.
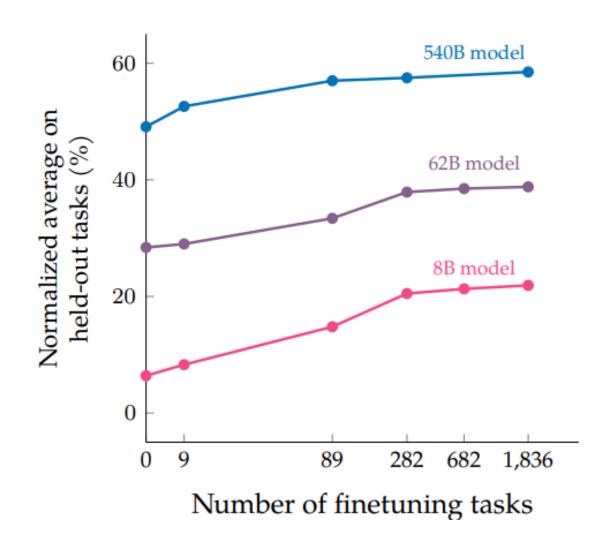
- Models' Sizes :

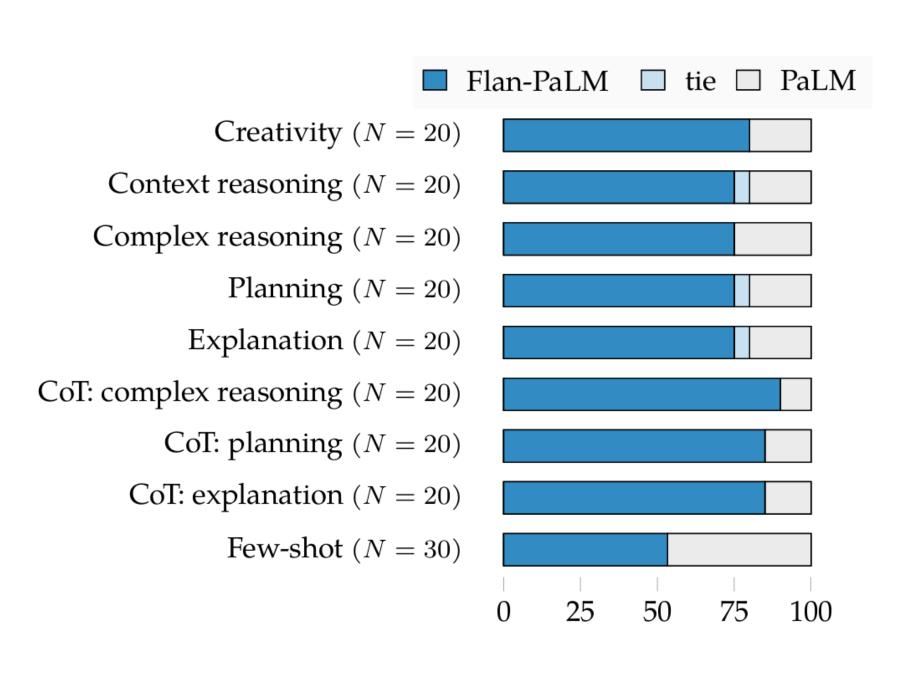| Params | Model | Architecture | Pre-training Objective |
|--------|-------|--------------|------------------------|
| 80M | Flan-T5-Small | encoder-decoder | span corruption |
| 250M | Flan-T5-Base | encoder-decoder | span corruption |
| 780M | Flan-T5-Large | encoder-decoder | span corruption |
| 3B | Flan-T5-XL | encoder-decoder | span corruption |
| 11B | Flan-T5-XXL | encoder-decoder | span corruption |
| 8B | Flan-PaLM | decoder-only | causal LM |
| 62B | Flan-PaLM | decoder-only | causal LM |
| 540B | Flan-PaLM | decoder-only | causal LM |
| 62B | Flan-cont-PaLM | decoder-only | causal LM |
| 540B | Flan-U-PaLM | decoder-only | prefix LM + span corruption |

- Results:

**Performance improves on :**
  - **different model classes**
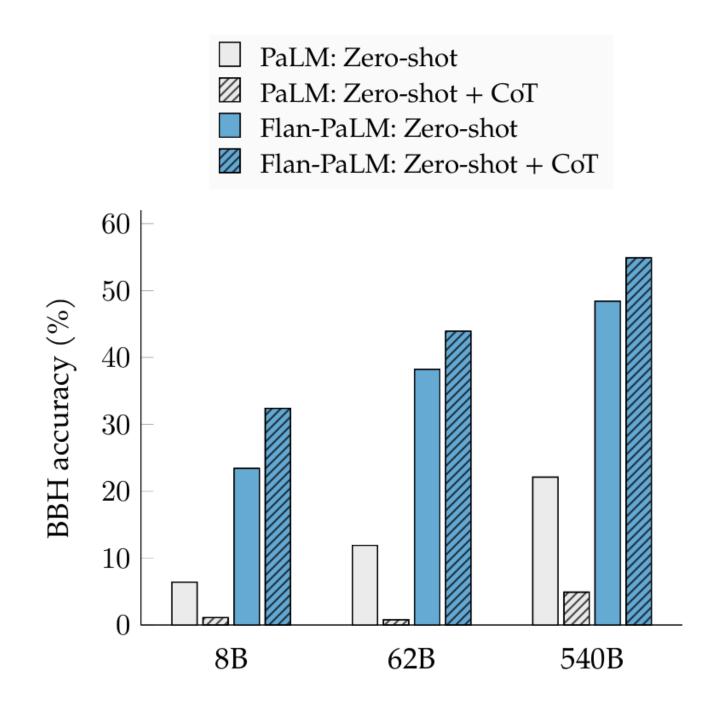  - **prompting setups**
  - **evaluation benchmarks**



| Params | Model | Norm. avg. | MMLU Direct | MMLU CoT | BBH Direct | BBH CoT | TyDiQA Direct | MGSM CoT |
|--------|-------|-----------|-------------|----------|------------|---------|---------------|----------|
| 80M | T5-Small | -9.2 | 26.7 | 5.6 | 27.0 | 7.2 | 0.0 | 0.4 |
|  | Flan-T5-Small | -3.1 (+6.1) | 28.7 | 12.1 | 29.1 | 19.2 | 1.1 | 0.2 |
| 250M | T5-Base | -5.1 | 25.7 | 14.5 | 27.8 | 14.6 | 0.0 | 0.5 |
|  | Flan-T5-Base | 6.5 (+11.6) | 35.9 | 33.7 | 31.3 | 27.9 | 4.1 | 0.4 |
| 780M | T5-Large | -5.0 | 25.1 | 15.0 | 27.7 | 16.1 | 0.0 | 0.3 |
|  | Flan-T5-Large | 13.8 (+18.8) | 45.1 | 40.5 | 37.5 | 31.5 | 12.3 | 0.7 |
| 3B | T5-XL | -4.1 | 25.7 | 14.5 | 27.4 | 19.2 | 0.0 | 0.8 |
|  | Flan-T5-XL | 19.1 (+23.2) | 52.4 | 45.5 | 41.0 | 35.2 | 16.6 | 1.9 |
| 11B | T5-XXL | -2.9 | 25.9 | 18.7 | 29.5 | 19.3 | 0.0 | 1.0 |
|  | Flan-T5-XXL | 23.7 (+26.6) | 55.1 | 48.6 | 45.3 | 41.4 | 19.0 | 4.9 |
| 8B | PaLM | 6.4 | 24.3 | 24.1 | 30.8 | 30.1 | 25.0 | 3.4 |
|  | Flan-PaLM | 21.9 (+15.5) | 49.3 | 41.3 | 36.4 | 31.1 | 47.5 | 8.2 |
| 62B | PaLM | 28.4 | 55.1 | 49.0 | 37.4 | 43.0 | 40.5 | 18.2 |
|  | Flan-PaLM | 38.8 (+10.4) | 59.6 | 56.9 | 47.5 | 44.9 | 58.7 | 28.5 |
| 540B | PaLM | 49.1 | 71.3 | 62.9 | 49.1 | 63.7 | 52.9 | 45.9 |
|  | Flan-PaLM | 58.4 (+9.3) | 73.5 | 70.9 | 57.9 | 66.3 | 67.8 | 57.0 |
| 62B | cont-PaLM | 38.1 | 61.2 | 57.6 | 41.7 | 53.1 | 45.7 | 32.0 |
|  | Flan-cont-PaLM | 46.7 (+8.6) | 66.1 | 62.0 | 51.0 | 53.3 | 62.7 | 40.3 |
| 540B | U-PaLM | 50.2 | 71.5 | 64.0 | 49.2 | 62.4 | 54.6 | 49.9 |
|  | Flan-U-PaLM | 59.1 (+8.9) | 74.1 | 69.8 | 59.3 | 64.9 | 68.3 | 60.4 |

# Flan-T5 : Scaling Instruction-Finetuned Language Models

- Results:



**Instruction Tuning Improves Human Usability**



**Finetuning with Chain-of-Thought Data Unlocks Zero-Shot Reasoning**

# References

- Chung, H. W., Hou, L., Longpre, S., Zoph, B., Tay, Y., Fedus, W., Li, Y., Wang, X., Dehghani, M., Brahma, S., Webson, A., Gu, S. S., Dai, Z., Suzgun, M., Chen, X., Chowdhery, A., Castro-Ros, A., Pellat, M., Robinson, K.,. . . Wei, J. (2022, 20 octobre). Scaling Instruction-Finetuned Language Models. arXiv.org. https://arxiv.org/abs/2210.11416

- Sanh, V., Webson, A., Raffel, C., Bach, S. H., Sutawika, L., Alyafeai, Z., Chaffin, A., Stiegler, A., Scao, T. L., Raja, A., Dey, M., Bari, M. S., Xu, C., Thakker, U., Sharma, S. S., Szczechla, E., Kim, T., Chhablani, G., Nayak, N.,. . . Rush, A. M. (2021, 15 octobre). Multitask prompted training enables Zero-Shot task generalization. arXiv.org. https://arxiv.org/abs/2110.08207
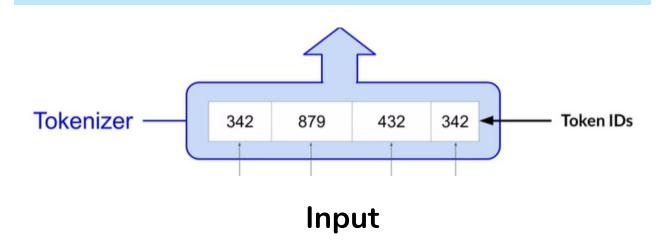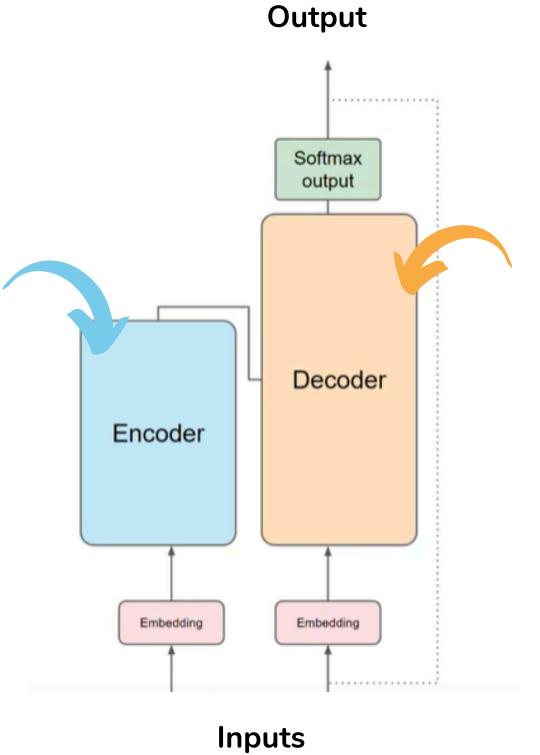
THANK YOU

CentraleSupélec

**sequence vectors:**

encodes information about both the element itself and its context within the sequence.

**process the input sequence and create contextualized representations for each element (word or token) in the sequence**

Output

Softmax output

Decoder

Encoder

Embedding    Embedding

Inputs

**the final representation vectors**

**generate an output sequence based on the contextualized representations provided by the encoder**

Tokenizer ——— | 342 | 879 | 432 | 342 | ◄——— Token IDs

**Input**