

Instruction Tuning with FLAN

Finetuned Language Models are Zero-Shot Learners

Google Research



Jason Wei



Maarten Bosma



Vincent Zhao



Kelvin Guu



Adams Yu



Brian Lester



Nan Du



Andrew Dai



Quoc Le



Summary

- “**Instruction tuning**” finetunes a language model on a collection of NLP tasks described using instructions.
- We instruction-tune a 137B pretrained LM and call the resulting model **FLAN** (for Finetuned Language Net).
- Instruction tuning helps the model perform tasks it wasn’t trained on, giving the model a range of applications.

Sentiment Analysis

Input:

- Movie review: This movie is the best RomCom since Pretty Woman.

Did this critic like the movie?

OPTIONS

- yes
- no

FLAN output:

- yes

Sentiment Analysis

Input:

- Movie review: This movie is the best RomCom since Pretty Woman.

Did this critic **dislike** the movie?

OPTIONS

- yes
- no

FLAN output:

- no

Question Answering

Input:

- What is Vincent Zhao's profession?

FLAN output:

- actor
- martial artist
- kung fu actor and martial artist

Question Answering

Input:

- What is Vincent Zhao's profession?

Answer in Chinese.

FLAN output:

- 演员
- 演員
- 导演

Rule-based tasks

Input:

- Jason Wei is reading the paper "Finetuned Language Models are Zero-Shot Learners."

Change to past tense.

FLAN output:

- Jason Wei was reading the paper "Finetuned Language Models are Zero-Shot Learners."

Rule-based tasks

Input:

- Jason Wei is reading the paper "Finetuned Language Models are Zero-Shot Learners."

Change to the verb to eat.

FLAN output:

- Jason Wei is eating the paper "Finetuned Language Models are Zero-Shot Learners."

Rule-based tasks

Input:

- Jason Wei is reading the paper "Finetuned Language Models are Zero-Shot Learners."

Change to passive voice.

FLAN output:

- The paper "Finetuned Language Models are Zero-Shot Learners" is being read by Jason Wei.

Rule-based tasks

Input:

- Jason Wei is reading the paper "Finetuned Language Models are Zero-Shot Learners."

Change the subject to Dipanjan Das and the paper to "BERT Rediscovered the NLP Pipeline".

FLAN output:

- Dipanjan Das is reading the paper "BERT rediscovered the NLP pipeline".

Recommendations

Input:

- Recommend activities to do on a **sunny weekend in Mountain View.**

FLAN output:

- Visit the Computer History Museum
- Hiking
- Picnic
- Go to the beach
- Visit Winchester Mystery House
- Wineries
- Bike ride

FLAN can generate training data for downstream classifiers.
Maybe you want to classify new intents.

Input:

- Generate utterances with the intent "**get COVID vaccine**"

FLAN output:

- I would like to get the covid vaccine.
- I would like a COVID vaccine as soon as possible
- Can I get a COVID vaccine as a preventative measure?
- get COVID vaccine for a five year old
- where is a center for COVID vaccine testing?

Style transfer

Input:

- "I guess I need to do laundry."

Rewrite the above sentence in a Shakespearean style.

FLAN output:

- To my mind, verily, I suppose 'tis time, I wash my clothes.
- Mine eyes and ears and tongue declare me idle; and yet it is true I guess, By my garments, that I need to do laundry
- Lo my duty now calls To thee, to wash this body's clothes.

Or maybe you're like James Joyce and you like to make up words.

Input:

- Make up a new word that means "when bananas shower from the sky"

FLAN output (cherrypicked):

- Bananation
- Banana shower
- Banania
- Bananafall
- Bananada

Or just have FLAN write the story for you :)

Input:

- Write a sad story about carrots.

FLAN output:

- **The sun had barely risen when the carrots were taken from their beds and brought to market.** They'd had no time to gather their strength. They smelled of earth, their greens still wet from the dew. **At first no one was interested in buying them. But one man finally reached out and picked up one of the carrots. "I'll take this one!" he cried out. "It's for my mother! I'll pay you well for it."** The man took the carrot to the cashier, eager to get it home to his mother. He was careful to wrap it carefully so that it wouldn't get bruised or dirty. **He hurried home, and when he saw his mother, he gave her the carrot. "Is this one of the ones I grew?" she asked, her voice wavering with emotion.** "Yes, mother," he replied. "It is, dear." The man left the carrot with his mother and hurried back to market. **His mother was very sick, and he was afraid she wouldn't last much longer. He needed to buy all the carrots in market so that none would be left when she was gone.**

Failure cases

Input:

- Jason Wei is reading the paper "Finetuned Language Models are Zero-Shot Learners."

What is the second word in the sentence?

FLAN output:

- paper
- Jason Wei
- learners
- language

Input:

- What is Vincent Zhao's profession?

Answer in Danish.

FLAN output:

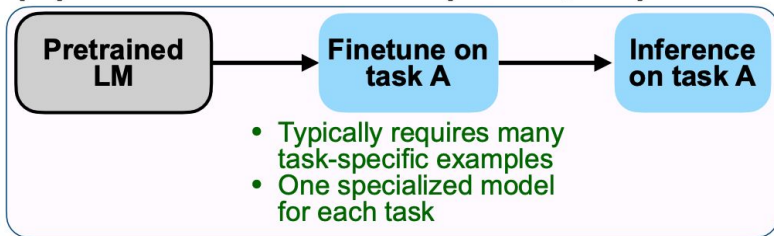
- Hvad er Vincent Zhao's erhverv?
- Hvad er Vincent Zhao's beskæftigelse?

Talk Outline

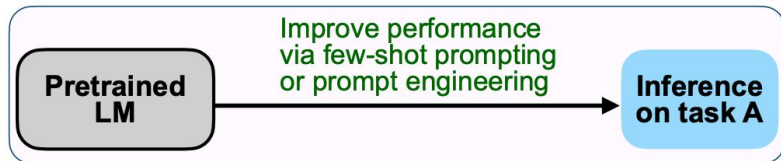
1. Background and motivation
2. Training FLAN & experimental setup
3. Results on various tasks, ablation studies

Motivation

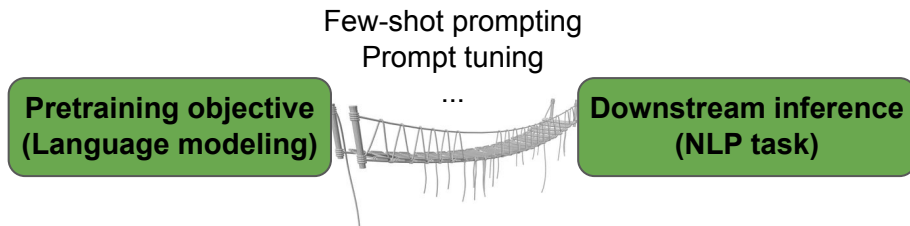
(A) Pretrain–finetune (BERT, T5)



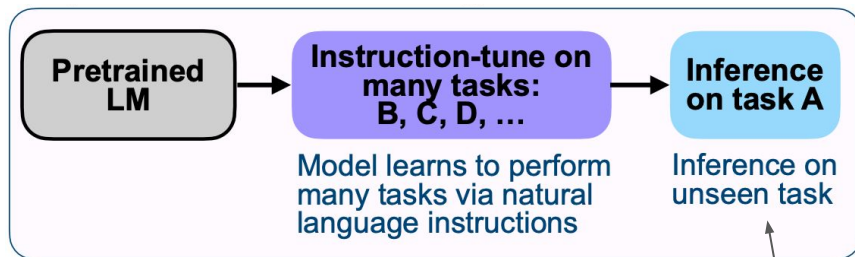
(B) Prompting (GPT-3)



`"This movie sucks." This movie review is {negative,positive}.`



Instruction tuning



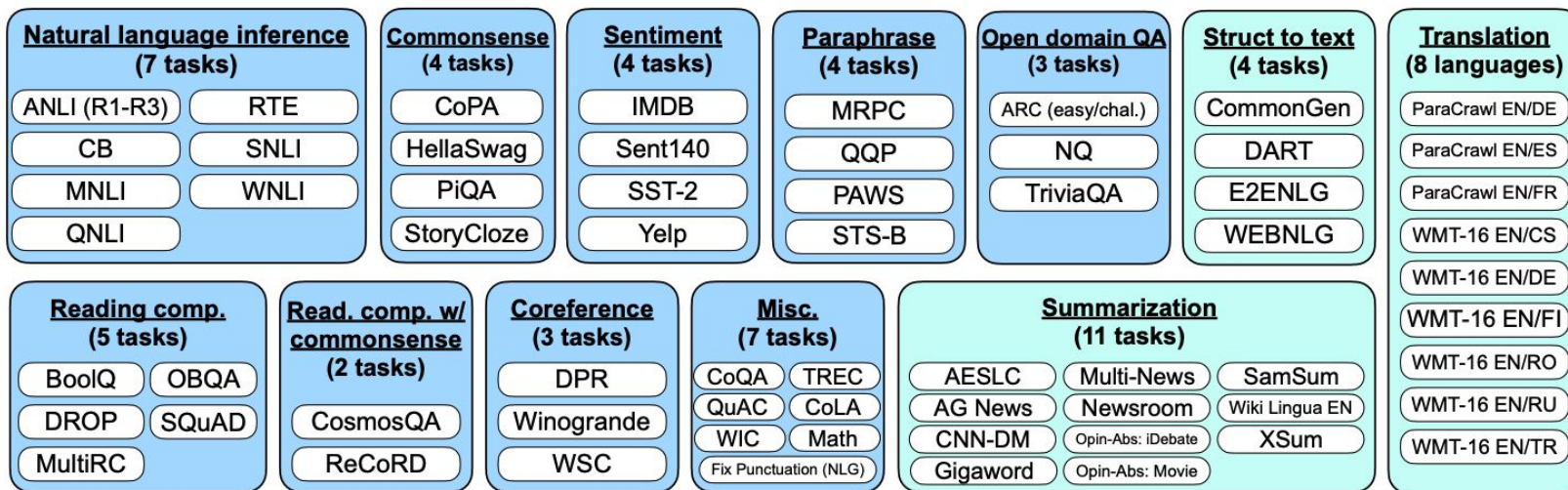
Can we use “a little bit” of supervision to teach the model to perform many NLP tasks?

i.e., zero-shot!

“Instruction tuning”—finetuning a language model on a collection of tasks described via instructions—improves the zero-shot performance of language models on unseen tasks.

NLP tasks and datasets

- 62 NLP datasets
- 12 “task clusters”



Templates

Premise

Russian cosmonaut Valery Polyakov set the record for the longest continuous amount of time spent in space, a staggering 438 days, between 1994 and 1995.

Hypothesis

Russians hold the record for the longest stay in space.

Target

Entailment
Not entailment



Options:
- yes
- no



Template 1

<premise>

Based on the paragraph above, can we conclude that <hypothesis>?

<options>

Template 2

<premise>

Can we infer the following?

<hypothesis>

<options>

Template 3

Read the following and determine if the hypothesis can be inferred from the premise:

Premise: <premise>

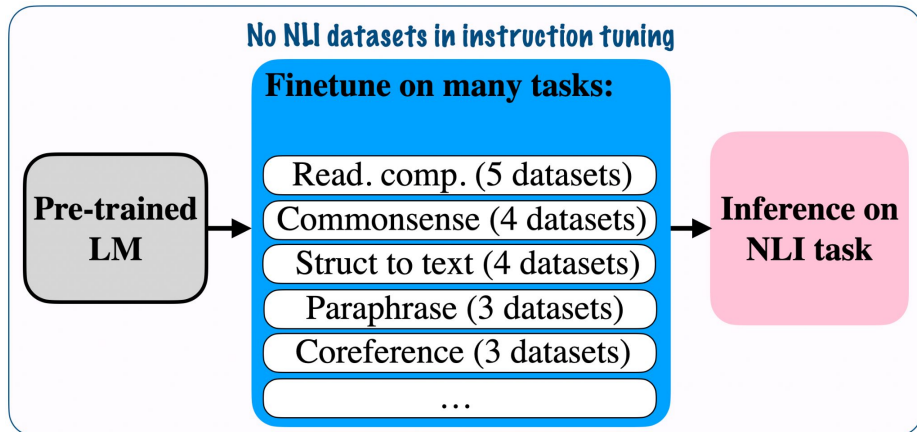
Hypothesis: <hypothesis>

<options>

Template 4. ...

We generate many natural instruction templates for each task

Evaluation splits



We evaluate on “unseen” / “zero-shot” tasks where no datasets from that task were seen during instruction tuning.

Classification with “options”

Input (Commonsense Reasoning)

Here is a goal: Get a cool sleep on summer days.

How would you accomplish this goal?

OPTIONS:

-Keep stack of pillow cases in fridge.

-Keep stack of pillow cases in oven.

Target

keep stack of pillow cases in fridge

For classification tasks, we teach FLAN to return one of several “options”

FLAN Training details

- 137B parameter pretrained checkpoint
- Instruction tune for 30k steps on 62 datasets spanning 12 task clusters

Summary of results

- 25 datasets spanning NLI, reading comprehension, closed-book QA, commonsense reasoning, coreference resolution, and translation
- Baselines: Base LM, GPT-3 175B



FLAN almost always outperforms Base LM

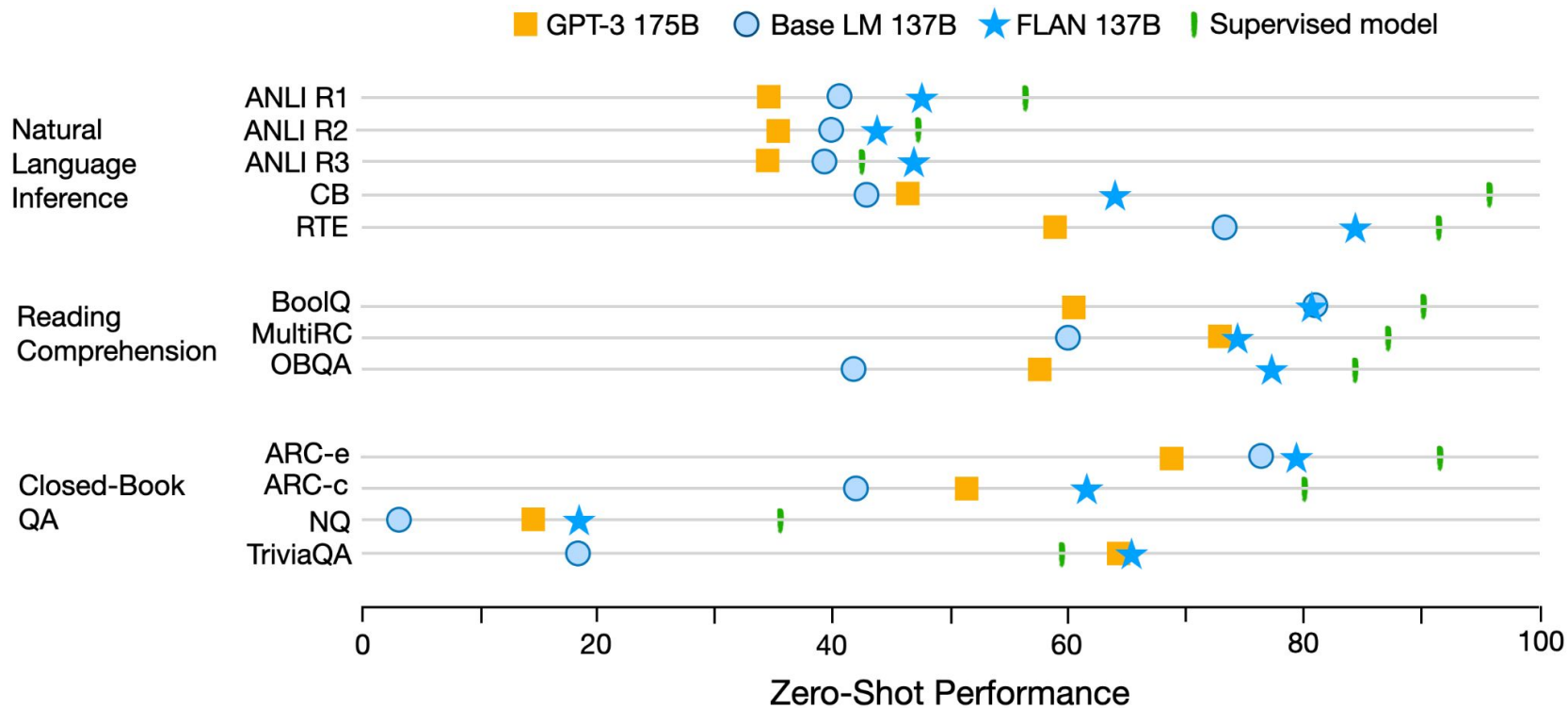


On 20 of 25 tasks, zero-shot FLAN outperforms zero-shot GPT-3

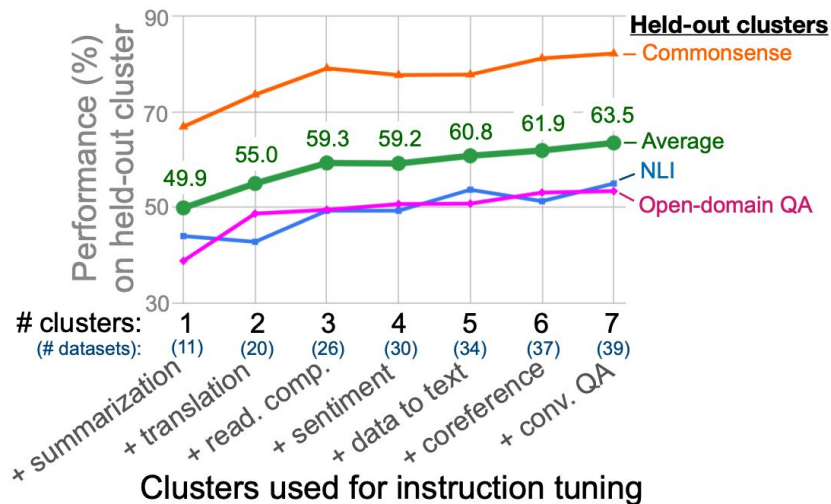


On 10 tasks, zero-shot FLAN even outperforms few-shot GPT-3

Results: NLI, reading comprehension, closed-book QA

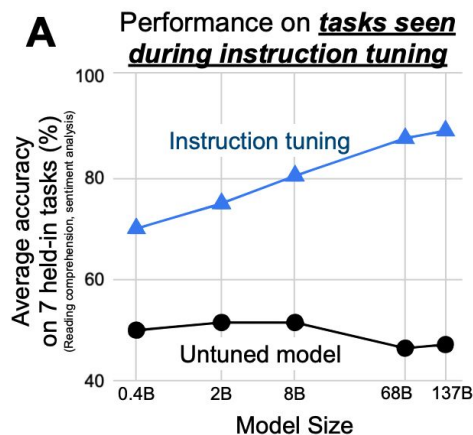


Ablation study: number of instruction tuning clusters

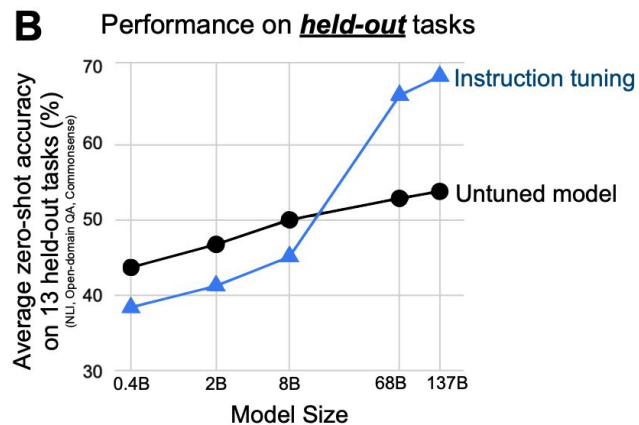


Adding additional task clusters to instruction tuning improves zero-shot performance on held-out task clusters.

Ablation study: scaling laws

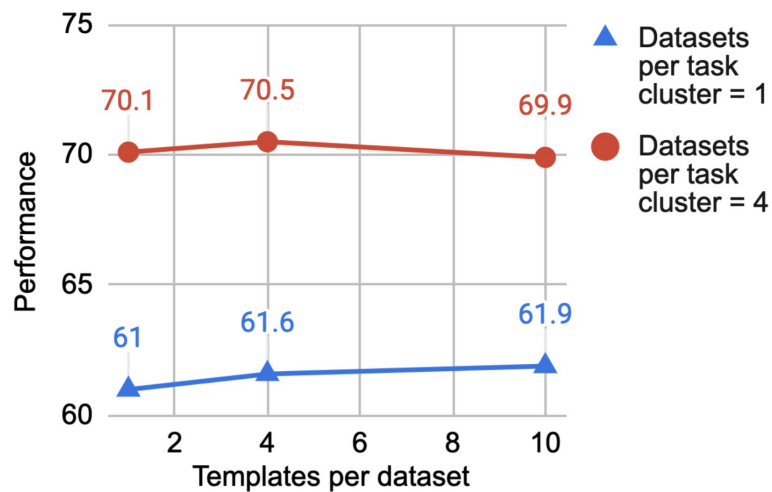


As expected, instruction tuning improves performance on seen tasks



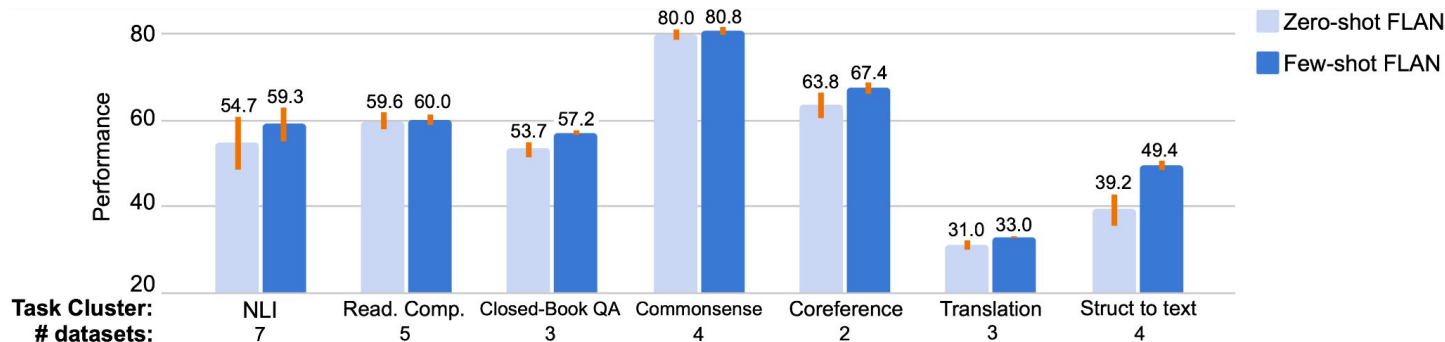
Performance on unseen tasks, on the other hand, only improves with sufficient model scale.

Ablation: templates per task



Curiously, more templates per dataset did not help much.

Further analysis: few-shot prompting



Few-shot prompting is a complementary way of improving performance with instruction tuning.

Example of few-shot prompt:

```
Does the following review have a positive or negative opinion of the movie?
```

```
<review>
```

```
Negative.
```

```
Does the following review have a positive or negative opinion of the movie?
```

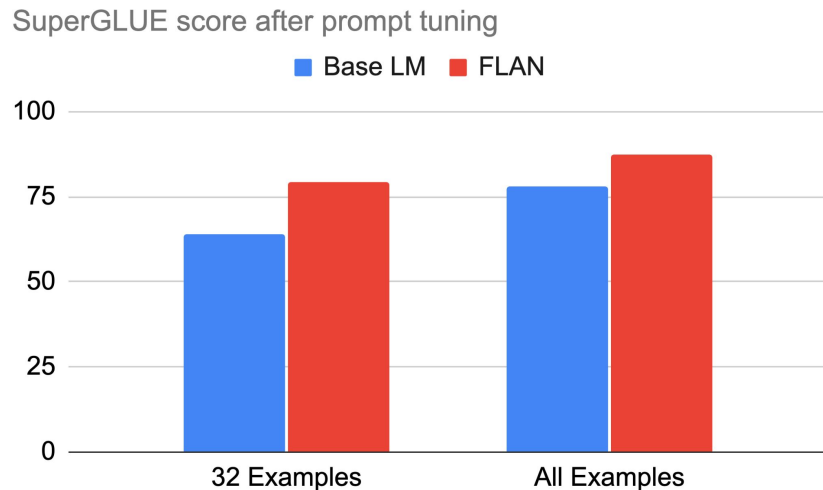
```
<review>
```

```
Positive.
```

```
Does the following review have a positive or negative opinion of the movie?
```

```
<review>
```

Further analysis: prompt tuning



FLAN responds better to continuous inputs from prompt tuning than base LM.

MULTITASK PROMPTED TRAINING ENABLES ZERO-SHOT TASK GENERALIZATION

Victor Sanh*
Hugging Face

Albert Webson*
Brown University

Colin Raffel*
Hugging Face

Stephen H. Bach*
Brown University

Lintang Sutawika
RioScience

Zaid Alyafeai
KFIIPM

Antoine Chaffin
IRISA IMATAG

Arnaud Stiegler
Hyperscience

Arun Raja
A*STAR

Colin Raffel/BigScience
did something similar
recently, on Oct 15
2021

	GPT-3	T0	FLAN
Size	175B	11B	137B
Multitask supervision	Implicit	Explicit	Explicit
Architecture/pretraining	Decoder/LM	Encoder-decoder/MLM+LM	Decoder/LM
Prompts	N/A	170 datasets* 1,939 prompts	62 datasets 620 prompts
Performance	Intriguing	Only available for NLI, Story completion, coreference, and some of BIG-bench, but also intriguing	Better than GPT-3 on average

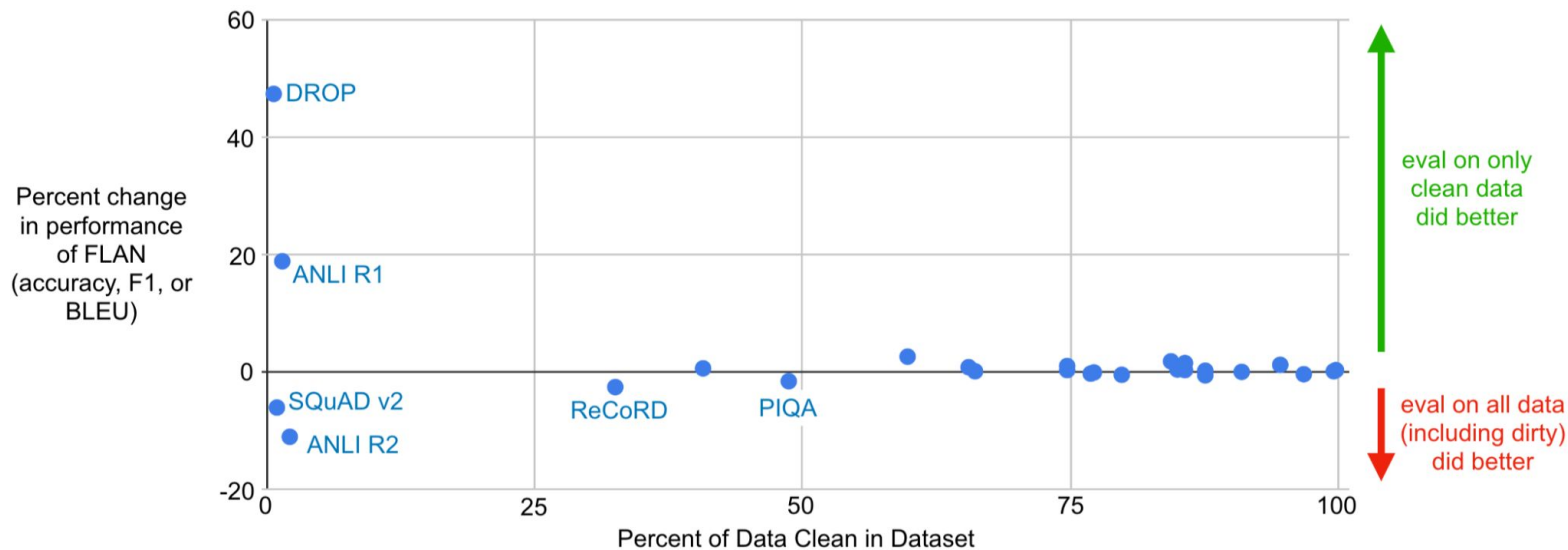
Conclusions

- Finetuning a language model on a collection of tasks allows it to follow instructions for a new task.
- This instruction-tuned language model has better zero-shot performance.
- Number of instruction tuning clusters and model scale are crucial.

Questions?

<u>Input</u>	<u>T0pp Output</u>
<p>Jason Wei is reading the paper "Finetuned Language Models are Zero-Shot Learners."</p> <p>Change to past tense.</p>	<p>Jason Wei is reading the paper "Finetuned Language Models are Zero-Shot Learners."</p>
<p>Jason Wei is reading the paper "Finetuned Language Models are Zero-Shot Learners."</p> <p>Change the verb to eat.</p>	<p>eat</p>
<p>Jason Wei is reading the paper "Finetuned Language Models are Zero-Shot Learners."</p> <p>Change to passive voice.</p>	<p>Jason Wei is reading the paper "Finetuned Language Models are Zero-Shot Learners."</p>
<p>Recommend activities to do on a sunny weekend in Mountain View.</p>	<p>Mountain View, California</p>
<p>Generate utterances with the intent "get COVID vaccine".</p>	<p>A nurse is giving a child a COVID vaccine.</p>

Further analysis: data contamination



We do not find evidence that example overlaps with pretraining data affects the performance of FLAN.