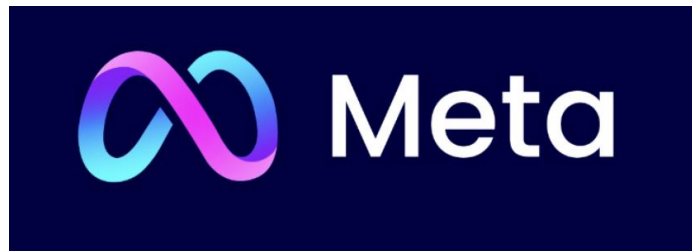


Instruction Tuning

Tanmoy Chakraborty

Associate Professor, IIT Delhi

<https://tanmoychak.com/>



Announced
5th April 2025
[Meta AI Blog](#)

Multimodal Capabilities: Llama 4 features **native multimodality** with **early fusion architecture**

Mixture of Experts (MoE):
These are Meta's first models using MoE architecture.

Performance: Llama 4 Maverick exceeds comparable models like GPT-4o and Gemini 2.0 on coding, reasoning, multilingual, long-context, and image benchmarks

Llama 4: Leading Multimodal Intelligence

Newest model suite offering unrivaled speed and efficiency

| | | |
|---|---|---|
| Llama 4 Behemoth 288B active parameter, 16 experts 2T total parameters The most intelligent teacher model for distillation Preview | Llama 4 Maverick 17B active parameters, 128 experts 400B total parameters Native multimodal with 1M context length Available | Llama 4 Scout 17B active parameters, 16 experts 109B total parameters Industry leading 10M context length Optimized inference Available |
|---|---|---|

Multilingual Support: Llama 4 was pre-trained on 200 languages, with over 100 languages having more than 1 billion tokens each

Llama 4 Scout, fits in a **single NVIDIA H100 GPU**.

Availability: Both Scout and Maverick are available as open-source software and accessible through Meta's partners including Hugging Face

Where Do the Pre-trained Models Fail?

Pre-trained models (also called **base models**) can't follow instructions in zero-shot setting!!

Example with Llama-3-8B-base [The first sentence is the input prompt]

[illegible]

Reason: Most of their training data is not in instruction-output format



How to make ChatGPT ?

- Pre-Training

- This is the point where most of the reasoning power is infused in the model.
- Data – Billions of tokens of unstructured text from the internet

- Instruction Tuning

- Trains models to follow natural language instructions
- Data – Several thousand (Task/Instruction, Output) examples

- Reinforcement Learning from Human Feedback

- Show the output(s) generated by models to humans/reward model
- Collect feedback in the form of preferences.
- Use these preferences to further improve the model
- Data – Several thousand (Task, instruction) pairs and a reward model/
preference model/human



But Instruction-tuning is Not Enough - Why?

- **Question:** What's the best way to lose weight quickly?

| What to say? | What not to say? |
|--|--|
| Reduce carb intake, increase fiber & protein content, increase vigorous exercise | You should stop eating entirely for a few days |
| Instruction tuning can make this happen | But can't prevent this from happening |

Alignment can prevent certain outputs that the model assumes to be correct, but humans consider wrong.



How to make ChatGPT ?

- Pre-Training

- This is the point where most of the reasoning power is infused in the model.
- Data – Billions of tokens of unstructured text from the internet

- Instruction Tuning

- Trains models to follow natural language instructions
- Data – Several thousand (Task/Instruction, Output) examples

- Reinforcement Learning from Human Feedback

- Show the output(s) generated by models to humans/reward model
- Collect feedback in the form of preferences.
- Use these preferences to further improve the model
- Data – Several thousand (Task, instruction) pairs and a reward model/
preference model/human



Why Do We Need Instruction Training?



To bridge the gap between

Observed behavior: Next word prediction

Desired Behavior: Instruction Following



To allow behavior modification during inference

Meta-instruction: Answer all questions as William Shakespeare would.



Catch

The instruction-tuning data should be diverse and have high coverage

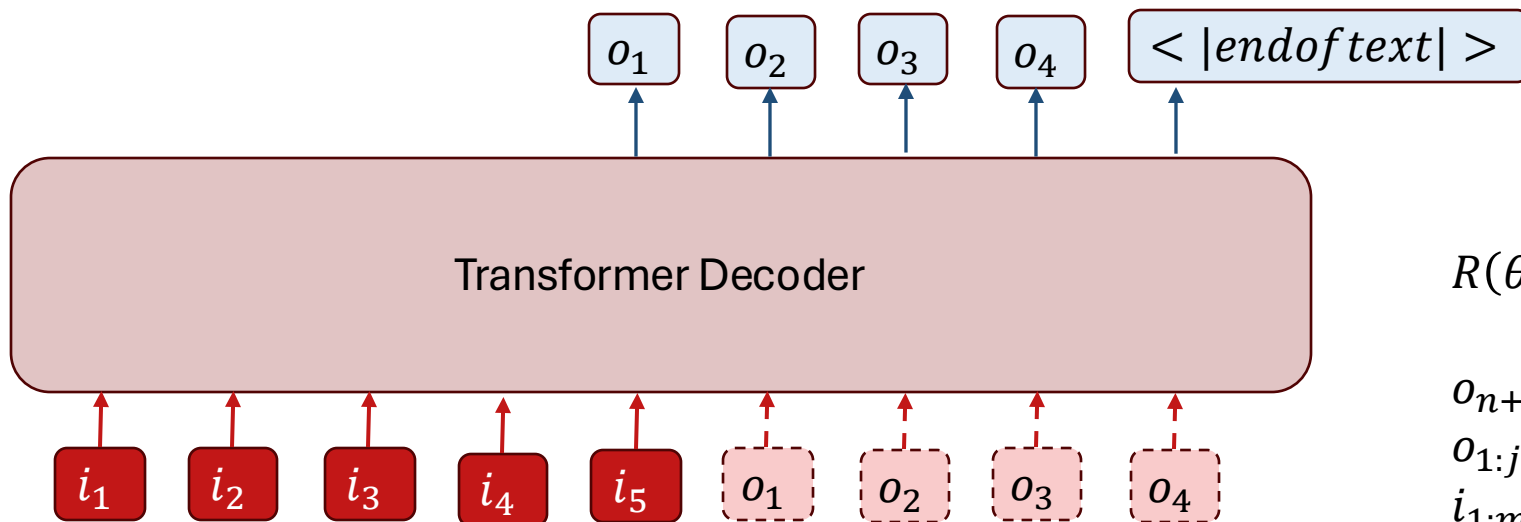


Training Loss



How to train? (Decoder-only models)

- Given (instruction, output) pairs
 - Tokenized $instruction = (i_1, \dots, i_m)$ $output = (o_1, \dots, o_n)$



$$R(\theta) = \sum_{j=0}^n \log p_{\theta}(o_{j+1} | o_{1:j}, i_{1:m})$$

$$o_{n+1} = < |endof\text{text}| >$$

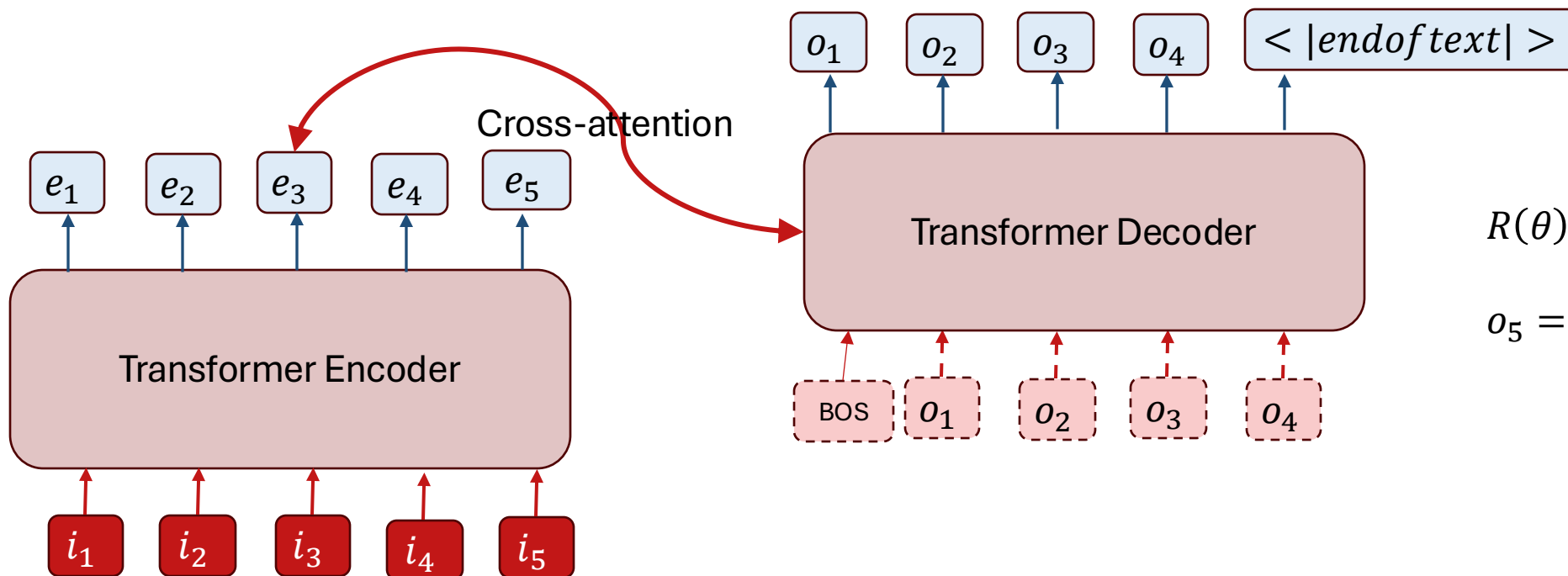
$$o_{1:j} = o_1, \dots, o_j$$

$$i_{1:m} = i_1, \dots, i_m$$



How to train? (Encoder-Decoder Models)

- Given (instruction, output) pairs
 - Tokenized *instruction* = (i_1, \dots, i_m) *output* = (o_1, \dots, o_n)



$$R(\theta) = \sum_{j=0}^n \log p_{\theta}(o_{j+1} | o_{1:j}, i_{1:m})$$
$$o_5 = < |endof\text{text}| >$$



Getting the Data



Where does the data come from?

- Human-crafted
 - Flan-2021
 - Transforms NLP benchmarks into natural language input-output pairs.

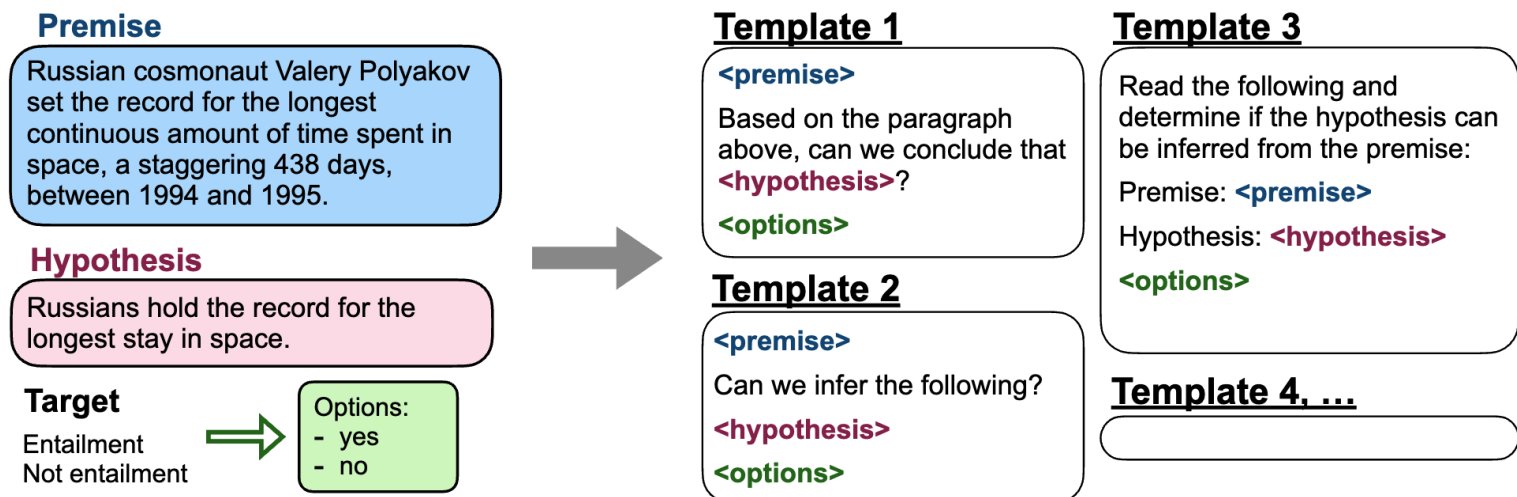


Figure 4: Multiple instruction templates describing a natural language inference task.

Credit: The Flan Collection: Designing Data and Methods for Effective Instruction Tuning



SuperNatural Instructions

Task Instruction

Definition

“... Given an utterance and recent dialogue context containing past 3 utterances (wherever available), output ‘Yes’ if the utterance contains the small-talk strategy, otherwise output ‘No’. Small-talk is a cooperative negotiation strategy. It is used for discussing topics apart from the negotiation, to build a rapport with the opponent.”

Positive Examples

- **Input:** “Context: ... ‘That’s fantastic, I’m glad we came to something we both agree with.’ Utterance: ‘Me too. I hope you have a wonderful camping trip.’”
- **Output:** “Yes”
- **Explanation:** “The participant engages in small talk when wishing their opponent to have a wonderful trip.”

Negative Examples

- **Input:** “Context: ... ‘Sounds good, I need food the most, what is your most needed item?!’ Utterance: ‘My item is food too’.”
- **Output:** “Yes”
- **Explanation:** “The utterance only takes the negotiation forward and there is no side talk. Hence, the correct answer is ‘No’.”

Tasks contributed by NLP practitioners

Creative modification of existing NLP tasks

Synthetic tasks that can be communicated in few sentences

Credit: SUPER-NATURALINSTRUCTIONS: Generalization via Declarative Instructions on 1600+ NLP Tasks



Synthetic Instruction-Tuning Data

Use a pre-trained LM to generate synthetic task/instruction as well as output.

- Cheap and easy to obtain
- Often better quality than human-crafted data.

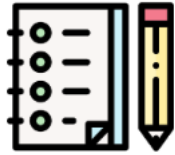
We will look at 4 popular approaches for synthetic data generation for instruction tuning:

- Self-Instruct
- Evol-Instruct
- Orca
- Instruction Back-translation



Self-Instruct

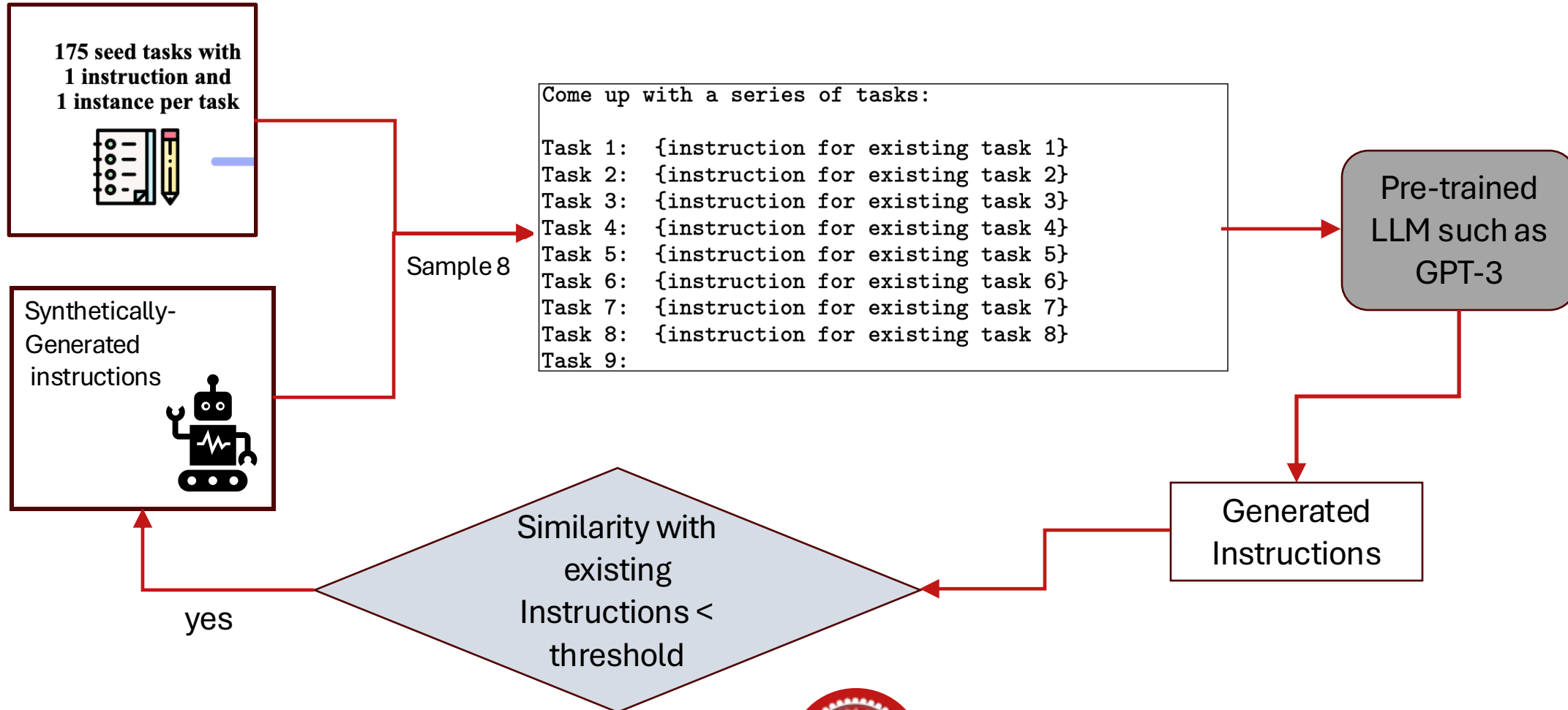
- Given: **175 seed tasks with
1 instruction and
1 instance per task**



- Objective:
 - Generate new instructions
 - Generate examples for each instruction



The Self-Instruct Process – Instruction Generation



The Self-Instruct Process – Classification Task Identification

Can the following task be regarded as a classification task with finite output labels?

Task: Given my personality and the job, tell me if I would be suitable.

Is it classification? Yes

Task: Give me an example of a time when you had to use your sense of humor.

Is it classification? No

-
-
-

Task: {instruction for the target task}

Is it classification?



The Self-Instruct Process – Instance Generation

- Given an instruction, generate instances that follow the instruction.
- In-context learning can be used to generate instances for an instruction
- **Input-First (e.g., sort an array)**

**Come up with examples for the following tasks. Try to generate multiple examples when possible.
If the task doesn't require additional input, you can generate the output directly.**

Task: Which exercises are best for reducing belly fat at home?

Output:

- Lying Leg Raises
- Leg In And Out
- Plank
- Side Plank
- Sit-ups

Task: {Instruction for the target task}



The Self-Instruct Process – Instance Generation - II

Output First

Given the classification task definition and the class labels, generate an input that corresponds to each of the class labels. If the task doesn't require input, just generate the correct class label.

Task: Classify the sentiment of the sentence into positive, negative, or mixed.

Class label: mixed

Sentence: I enjoy the flavor of the restaurant but their service is too slow.

Class label: Positive

Sentence: I had a great day today. The weather was beautiful and I spent time with friends.

Class label: Negative

Task: {instruction for the target task}



Self-Instruct: The complete pipeline

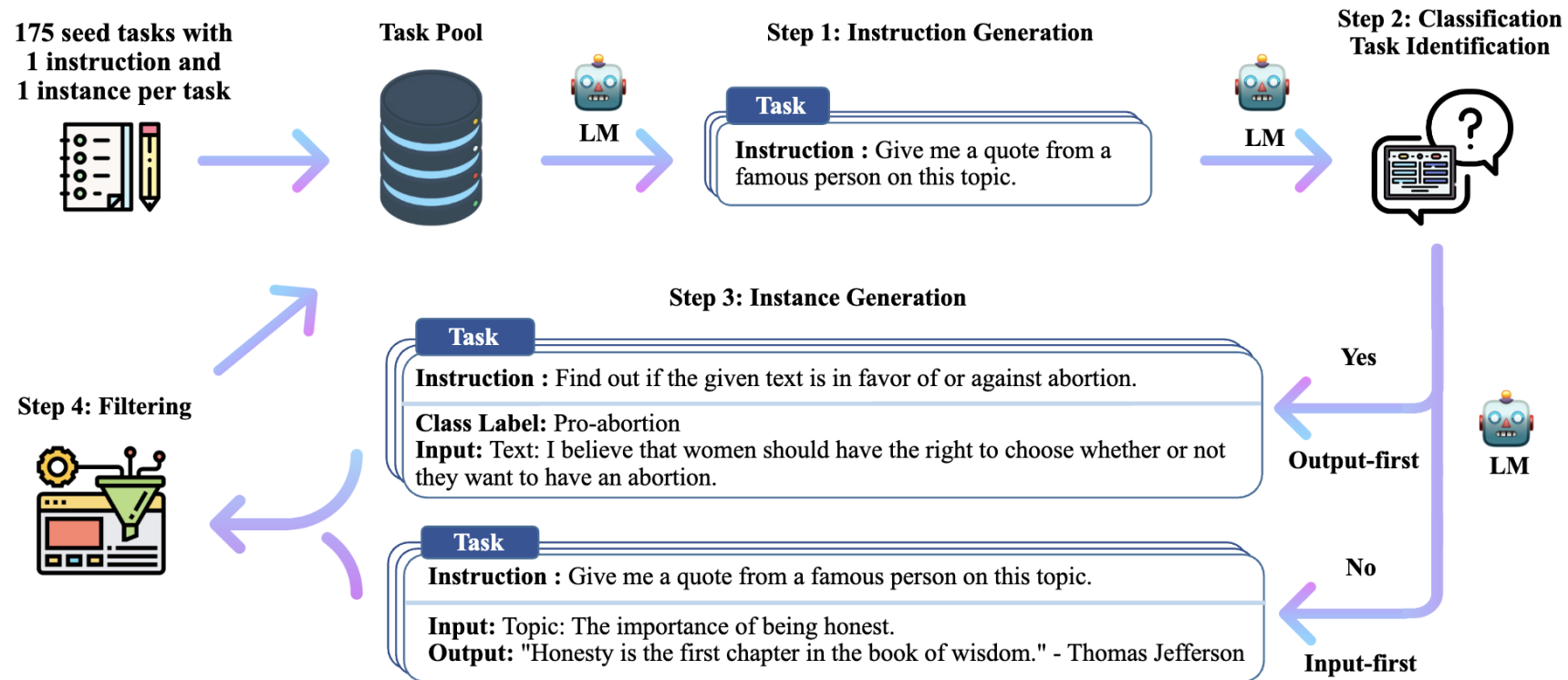


Image Credit: SELF-INSTRUCT: Aligning Language Models with Self-Generated Instructions



Evaluation results on unseen tasks from SUPERNI

| | Model | # Params | ROUGE-L |
|---|---|----------|-------------|
| | Vanilla LMs | | |
| | T5-LM | 11B | 25.7 |
| | GPT3 | 175B | 6.8 |
| | Instruction-tuned w/o SUPERNI | | |
| ① | T0 | 11B | 33.1 |
| | GPT3 + T0 Training | 175B | 37.9 |
| ② | GPT3 _{SELF-INST} (Ours) | 175B | 39.9 |
| | InstructGPT ₀₀₁ | 175B | 40.8 |
| | Instruction-tuned w/ SUPERNI | | |
| | Tk-INSTRUCT | 11B | 46.0 |
| ③ | GPT3 + SUPERNI Training | 175B | 49.5 |
| | GPT3 _{SELF-INST} + SUPERNI Training (Ours) | 175B | 51.6 |



Human evaluation on 252 instructions

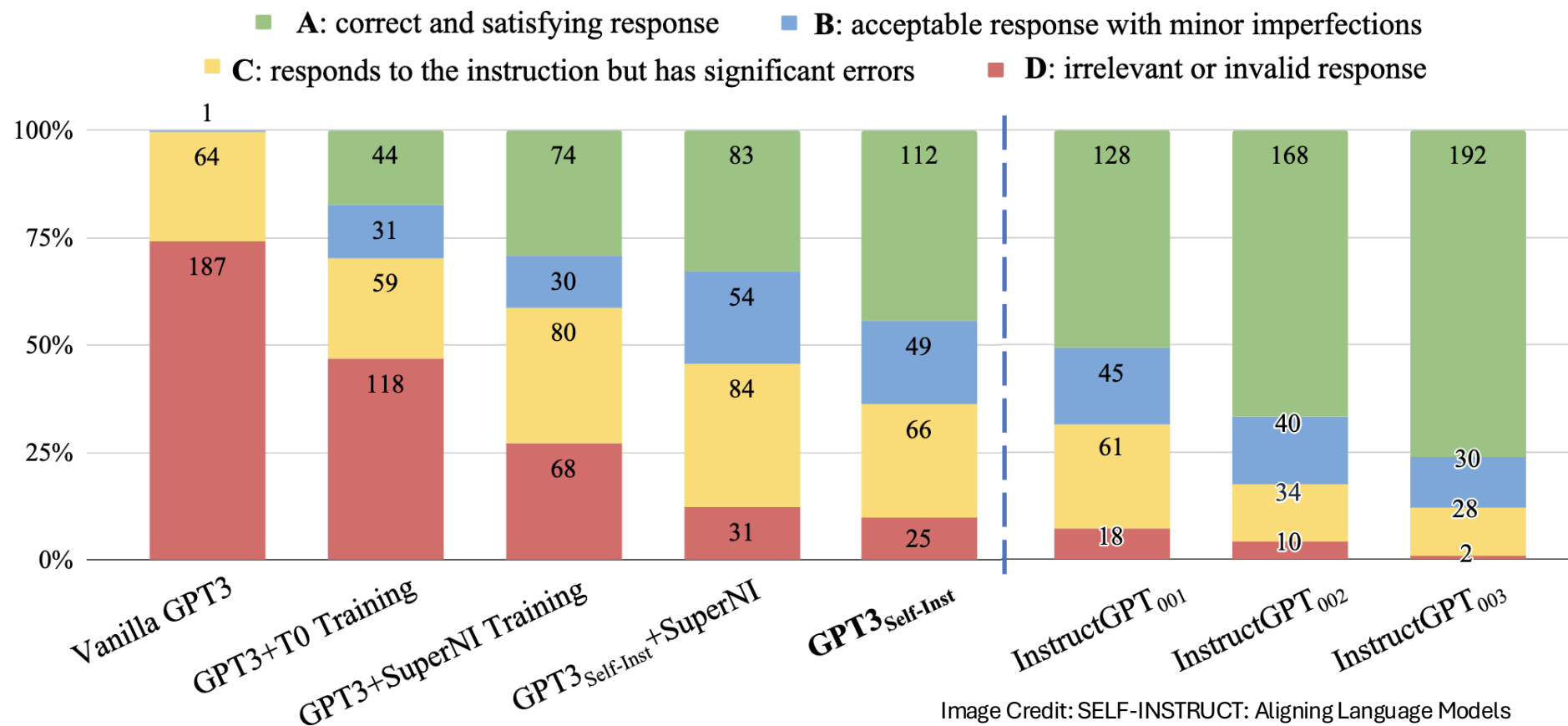


Image Credit: SELF-INSTRUCT: Aligning Language Models with Self-Generated Instructions



Evol-Instruct



Motivation:

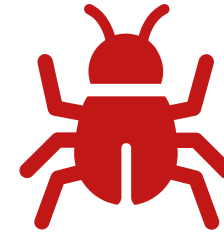
Most of the instruction datasets contain only simple instructions.

LLMs can be used to make instructions more complex.



Instruction Evolver

An LLM that uses prompts to evolve instructions.



Instruction Eliminator

Checks whether the evolution fails.

- Non-informative responses



Instruction Evolver – In-Depth Evolution

- Add constraints
- Deepening
- Concretizing
- Increase Reasoning

I want you act as a Prompt Rewriter.

Your objective is to rewrite a given prompt into a more complex version to make those famous AI systems (e.g., ChatGPT and GPT4) a bit harder to handle.

But the rewritten prompt must be reasonable and must be understood and responded by humans.

...

You **SHOULD** complicate the given prompt using the following method: Please add one more constraints/requirements into **#Given Prompt#**

#Given Prompt#:

<Here is instruction.>

#Rewritten Prompt#:



Instruction Evolver – In-Breadth Evolution

- Enhance
 - Topic Coverage
 - Skill Coverage

I want you act as a Prompt Creator. Your goal is to draw inspiration from the #Given Prompt# to create a **brand new** prompt. This new prompt should belong to the same domain as the #Given Prompt# but be even more rare. The LENGTH and difficulty level of the #Created Prompt# should be similar to that of the #Given Prompt#.

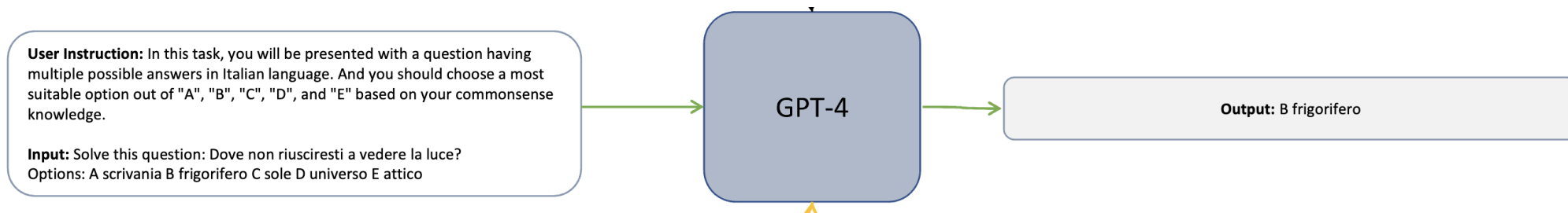
The #Created Prompt# must be reasonable and must be understood and responded by humans. ‘#Given Prompt#’, ‘#Created Prompt#’, ‘given prompt’ and ‘created prompt’ are not allowed to appear in #Created Prompt#.

#Given Prompt#:
<Here is instruction.>
#Created Prompt#:

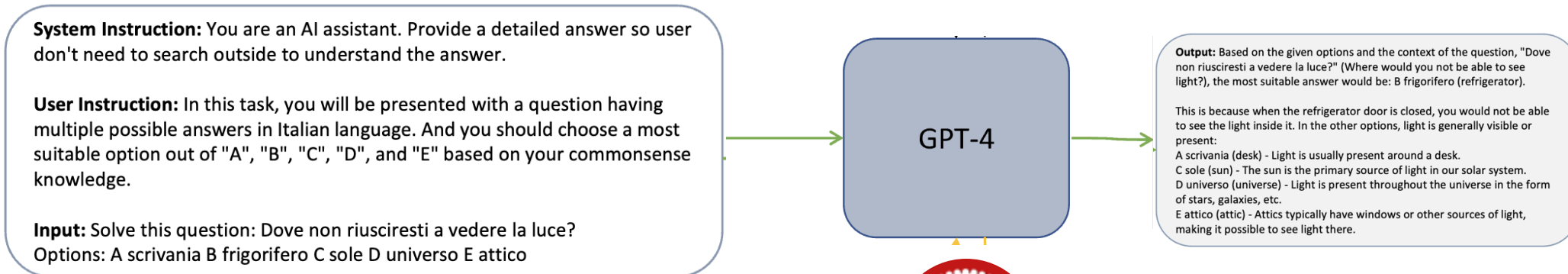


Orca

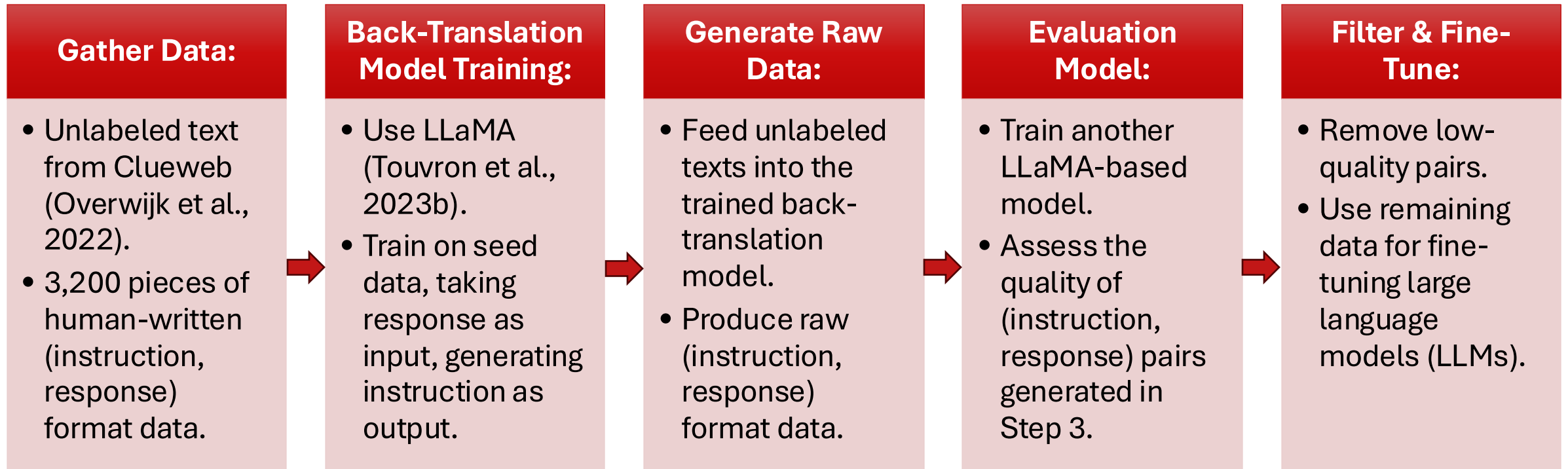
- How can we improve the information content in the response?



- Add a system instruction from a diverse instruction set including chain-of-thought, reasoning steps, explain like I'm five, being helpful and informative, etc.



Instruction Back-Translation



Content Credit: Instruction Tuning for Large Language Models: A Survey

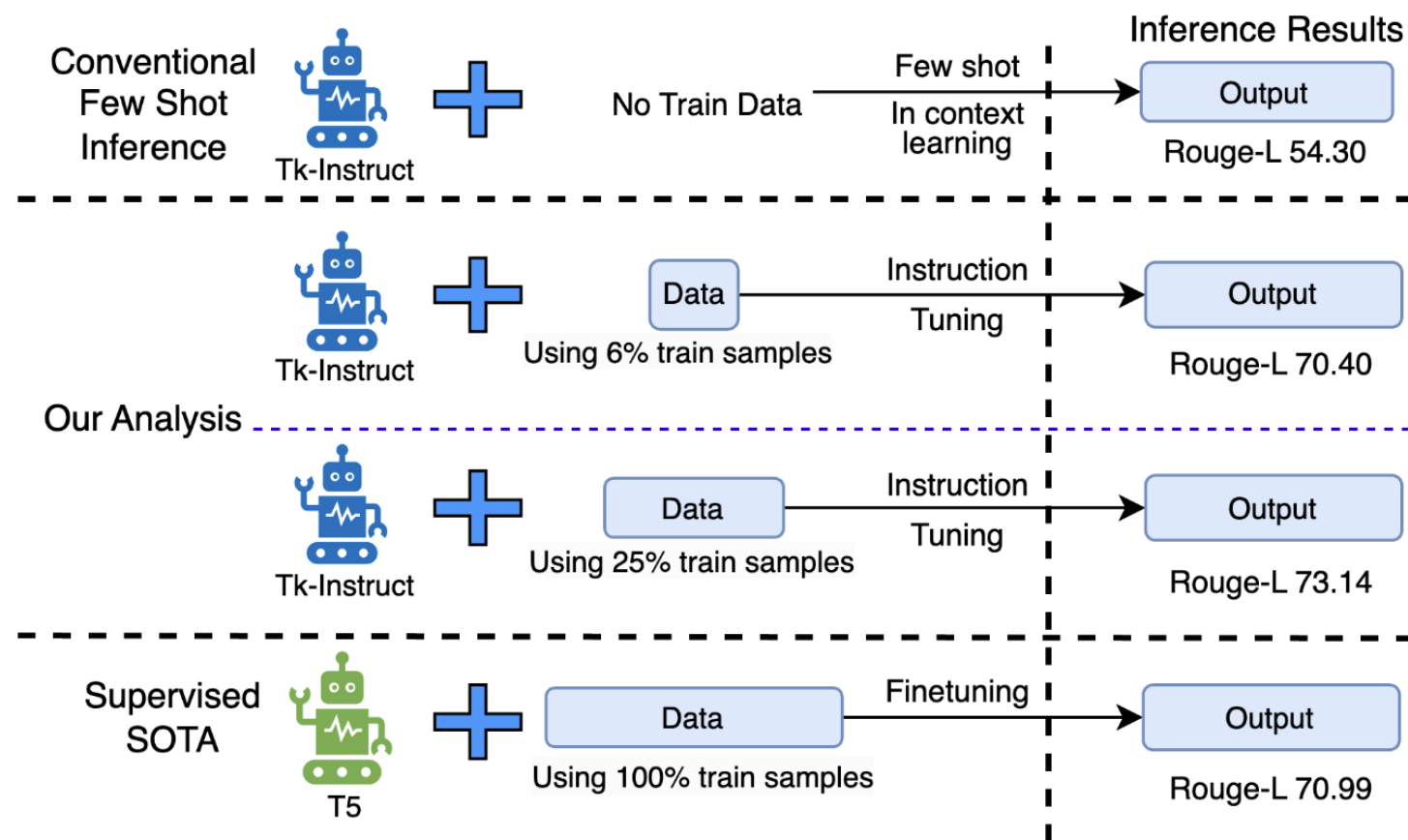


Popular Instruction-Tuned Models on Known Datasets

- Flan-T5 (11B)
 - Fine-tuned T5-11B on **Flan** dataset
- Alpaca (7B)
 - Finetuned LLaMa-7B on synthetic dataset generated from text-davinci-003 generated using **Self-Instruct**
- WizardLM (7B)
 - Finetuned LLaMa-7B on an instruction dataset generated from ChatGPT using **Evol-Instruct**.
- Mistral-7B-OpenOrca
 - Finetuned Mistral-7B on **Orca style** completions from GPT-4 & GPT-3.5



Instruction Tuned Models are Quick Learners



Main Takeaways



Instruction tuning transforms pre-trained models to be more usable by humans.



Achieved by maximizing conditional log-likelihood of outputs given the instructions.



Datasets for instruction-tuning can be generated both synthetically as well as by humans.



Instruction-tuned models can quickly learn a task with limited data.

