

# Automatic Reasoning and Tool-use (ART)

Combining CoT prompting and tools in an interleaved manner has shown to be a strong and robust approach to address many tasks with LLMs. These approaches typically require hand-crafting task-specific demonstrations and carefully scripted interleaving of model generations with tool use. [Paranjape et al., \(2023\)](#) propose a new framework that uses a frozen LLM to automatically generate intermediate reasoning steps as a program.

ART works as follows:

- given a new task, it select demonstrations of multi-step reasoning and tool use from a task library
- at test time, it pauses generation whenever external tools are called, and integrate their output before resuming generation

ART encourages the model to generalize from demonstrations to decompose a new task and use tools in appropriate places, in a zero-shot fashion. In addition, ART is extensible as it also enables humans to fix mistakes in the reasoning steps or add new tools by simply updating the task and tool libraries. The process is demonstrated below:

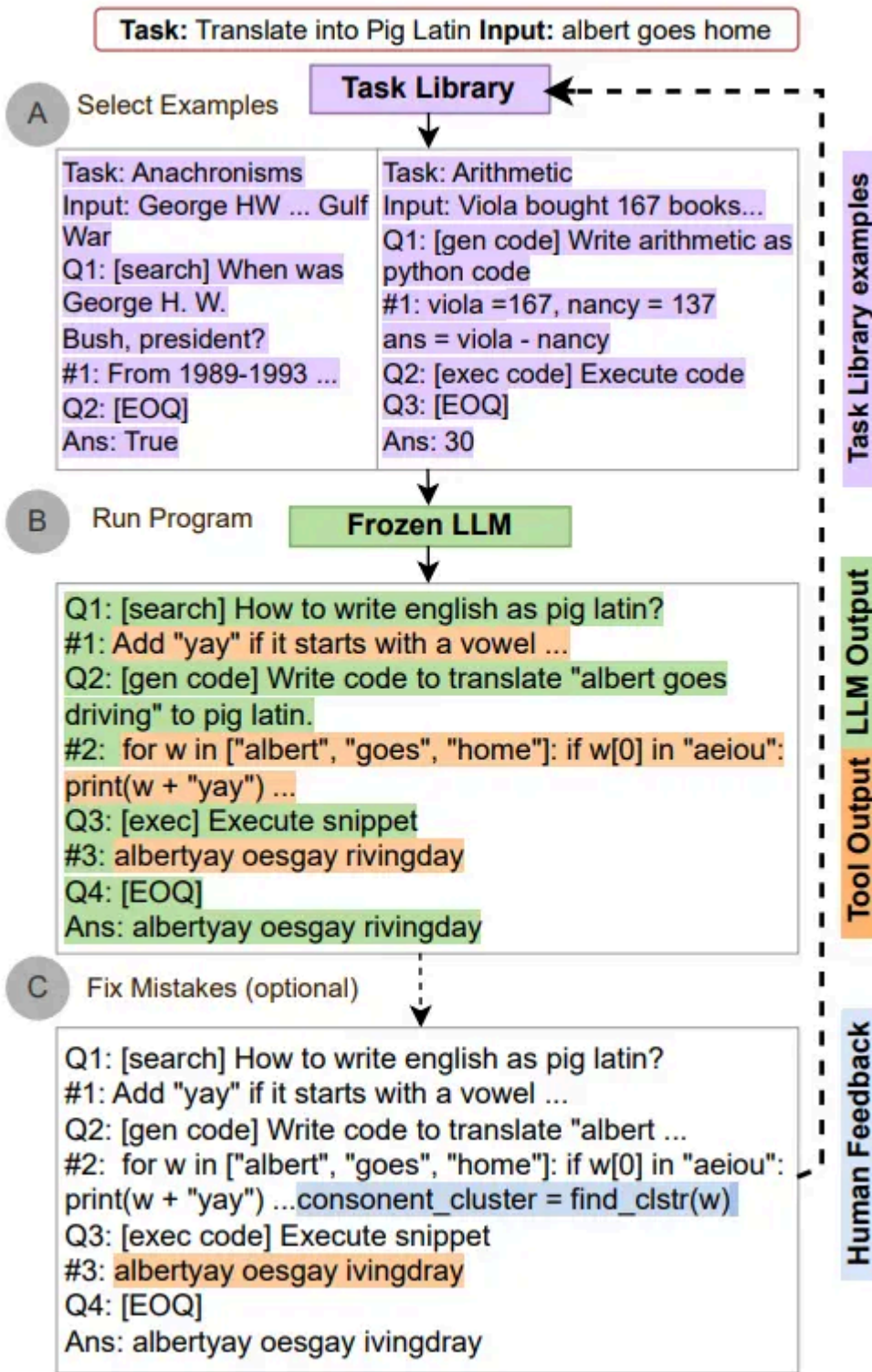


Image Source: [Paranjape et al., \(2023\)](#).

ART substantially improves over few-shot prompting and automatic CoT on unseen tasks in the BigBench and MMLU benchmarks, and exceeds performance of hand-crafted CoT prompts when human feedback is incorporated.

Below is a table demonstrating ART's performance on BigBench and MMLU tasks:

Task Name (Cluster)	Few Shot	AutoCot	ART w/o Tool Use	ART	GPT-3 Best
Test Tasks					
Sentence Ambiguity (Search)	70.67 <sup>5</sup>	51.47	71.00	73.33	-
Strategy QA (Search)	55.49 <sup>5</sup>	27.22	59.37	66.44	-
Physics (Search)	70.09 <sup>5</sup>	61.83	59.13	67.55	-
$\Delta$ with ART (Search)	<b>+3.7</b>	<b>+22.27</b>	<b>+ 5.9</b>		
Physics Questions (Arithmetic)	7.02 <sup>5</sup>	5.56	6.30	20.37	-
Operators (Arithmetic)	71.23 <sup>7</sup>	75.52	71.80	92.00	-
Unit interpretation (Arithmetic)	58.2 <sup>7</sup>	41.20	51.4	53.99	-
Repeat copy logic (Arithmetic)	50.01 <sup>7</sup>	15.63	31.25	44.38	-
Object Counting (Arithmetic)	39.2 <sup>7</sup>	26.80	42.2	87.00	81.20 <sup>1</sup>
Penguins in a table (Arithmetic)	58.23 <sup>7</sup>	40.40	68.86	77.85	72.34 <sup>1</sup>
Reasoning about objects (Arithmetic)	71.00 <sup>7</sup>	33.33	45.35	64.34	52.69 <sup>1</sup>
Tracking shuffled objects (Arithmetic)	22.39 <sup>7</sup>	19.44	18.14	37.67	36.32 <sup>1</sup>
$\Delta$ with ART (Arithmetic)	<b>+19.0</b>	<b>+36.7</b>	<b>+ 23.1</b>		<b>+6.1</b>
Word Unscramble (String)	40.72 <sup>7</sup>	32.44	23.03	42.7	-
Simple Text Editing (Code)	35.31 <sup>5</sup>	30.21	20.74	27.65	-
CS Algorithms (Code)	73.48 <sup>7</sup>	0.0	41.59	88.11	-
Sports Understanding (CoT)	69.74 <sup>5</sup>	51.47	92.89	-	86.59 <sup>1</sup>
Snarks (CoT)	54.58 <sup>5</sup>	57.24	57.13	-	65.2 <sup>1</sup>
Disambiguation QA (Free-form)	55.03 <sup>5</sup>	48.45	55.89	-	60.62 <sup>1</sup>
Temporal sequences (CoT)	55.80 <sup>7</sup>	19.70	49.5	-	81.8 <sup>1</sup>
Ruin names (CoT)	71.01 <sup>5</sup>	55.28	60.22	-	-
$\Delta$ with ART (Misc)	<b>2.4</b>	<b>22.5</b>	<b>24.37</b>		<b>-9.4</b>
$\Delta$ with ART (Overall)	<b>+6.9</b>	<b>+24.6</b>	<b>+16.7</b>		<b>-1.7</b>
MMLU					
College Computer Science (Search)	41.00	43.99	63.40	67.80	63.6 <sup>6</sup>
Astronomy (Search)	62.10	41.48	76.71	79.1	62.5 <sup>6</sup>
Business Ethics (Search)	61.60	48.8	77.17	81.16	72.7 <sup>6</sup>
Virology (Search)	50.03	49.52	71.60	71.49	50.72 <sup>6</sup>
Geography (Search)	77.67	57.07	70.30	71.71	81.8 <sup>6</sup>
Mathematics (Arithmetic)	36.67	33.77	39.50	45.66	34.5 <sup>6</sup>
$\Delta$ with ART (MMLU)	<b>+14.6</b>	<b>+23.7</b>	<b>+3.0</b>		<b>+8.5</b>

Table 3: ART performance on BigBench tasks and MMLU tasks. (<sup>1</sup> Human-crafted CoT (Wei et al., 2022; Suzgun et al., 2022), <sup>5</sup> InstructGPT (Ouyang et al., 2022), <sup>6</sup> Scaled instruction finetuning (Chung et al., 2022), <sup>7</sup> Code-davinci-002 (Chen et al., 2021)).

Image Source: [Paranjape et al., \(2023\)](#).

Last updated on September 19, 2024