# Language Pre-training without Natural Language

**Qian Liu** (刘乾) @ Sea AI Lab

**Bei Chen** (陈蓓) @ Microsoft Research Asia
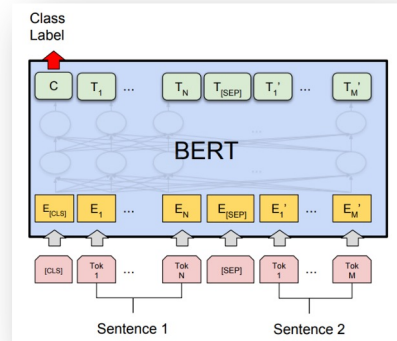
sigma

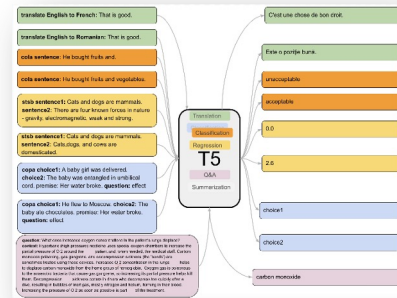# Current AI Paradigm: Language Models = SOTA





**BERT** (Devlin et al., 2018)



**T5** (Raffel et al., 2020)

# Current AI Paradigm: Language Models = Human Parity

# Research Challenge: **Reasoning**

However, the reasoning capability is still the mysterious for language models — even for giant language models (e.g., GPT3).

**Emergent Performance at 175B? No!**



**Reasoning, or correlation?**

# Research Challenge: **Reasoning**

However, it is difficult to obtain large amounts of clean natural language sentences containing clear evidence of reasoning.

# Key Idea: **Program as a Proxy**

There are rich reasoning operations (e.g., sort) in the program execution process. Can we leverage programs instead of natural language sentences as pre-training corpus?

| Program | | Natural Language |
|---|---|---|
| sorted([1, -5, 10, 6], key=abs, reverse=True) | $\approx$ | Given the list which contains 1, -5, 10 and 6, I want to order from high to low no matter what sign each number has, but keeping the sign |

# Key Idea: **Program as a Proxy**

There is a natural analogy between neural models and program executors!

| **Program** | **Natural Language** |
|---|---|
| sorted([1, -5, 10, 6], key=abs, reverse=True) | Given the list which contains 1, -5, 10 and 6, I want to order from high to low no matter what sign each number has, but keeping the sign |

**Program Executor**

**Neural Model**

[10, 6, -5, 1]

# Method Comparison: Execution v.s. Generation

Recent language models can perform program generation, and the difference is that we leverage program execution for natural language reasoning beyond programs.



*GitHub Copilot (2021)*

# Overview: Tabular, Numerical and Spatial Reasoning

# TAPEX: Table Pre-training via Learning a Neural SQL Executor

Qian Liu[1]  Bei Chen[2]  Jiaqi Guo[3]  Morteza Ziyadi[2]  Zeqi Lin[2]  Weizhu Chen[2]  Jian-Guang Lou[2]

1 北京航空航天大学 BEIHANG UNIVERSITY  2 Microsoft  3 西安交通大学 XI'AN JIAOTONG UNIVERSITY

# Background: **Tabular Reasoning**

| City | Country | Nations | Year |
|---|---|---|---|
| Athens | Greece | 14 | 1896 |
| St. Louis | USA | 12 | 1904 |
| ... | ... | ... | ... |
| Athens | Greece | 201 | 2004 |
| Beijing | China | 204 | 2008 |

## Question
Greece held its last Summer Olympics in which year

# Background: **Tabular Reasoning**

| City | Country | Nations | Year |
|---|---|---|---|
| Athens | Greece | 14 | 1896 |
| St. Louis | USA | 12 | 1904 |
| ... | ... | ... | ... |
| Athens | Greece | 201 | 2004 |
| Beijing | China | 204 | 2008 |

## Question

Greece held its last Summer Olympics in which year

# Previous Work: **Reinforcement Learning**

Obtain rewards by comparing execution results of sampled SQL queries with golden answers to train a text-to-SQL semantic parser. Hard to scale to complex scenarios.



[Chen et al. 2018]

# Previous Work: **Table Parsing**

Predict answer by selecting table cell values and optionally applying an aggregation operator to the selected region. Flexibility is limited.



[Herzig et al. 2020]

# Preliminary: Generative Language Model

We formulate the task of table-based question answering as answer generation, and leverages generative language models (e.g., BART) to output autoregressively.

ans1, ans2 … </s>

| Bidirectional Encoder | ⟹ | Autoregressive Decoder |

Question + Flattened Table

<s>, ans1, ans2 …

# Preliminary Result: Models Are Data-Hungry



WikiSQL (Weak)

TabFact

SQA

WikiTableQuestions

# Method: SQL Execution Pre-training

Pre-training a model to mimic the behavior of a symbolic execution engine.



Pre-trained LM for Textual Data → *Pre-training* → Pre-trained LM for Tabular Data → *Fine-tuning* → Fine-tuned LM for Table-related Task

**Input**: **SELECT** City **WHERE** Country = France **ORDER BY** Year **ASC LIMIT 1** [Table]
**Output**: Paris

Synthetic Pre-training Corpus

**Input**: Greece held its last Summer Olympics in which year? [Table]
**Output**: 2004

Realistic Downstream Datasets

#17

# Method: SQL Execution Pre-training

If we train a model to mimic the SQL query execution procedure over databases, we believe it learns programmatic reasoning from the execution engine.

take a table

sample an executable SQL query

| Year | City | Country | Nations |
|------|------|---------|---------|
| 1896 | Athens | Greece | 14 |
| 1900 | Paris | France | 24 |
| 1904 | St. Louis | USA | 12 |
| … | … | … | … |
| 2004 | Athens | Greece | 201 |
| 2008 | Beijing | China | 204 |

**SELECT** City **WHERE** Country = France **ORDER BY** Year **ASC LIMIT 1**

flatten

**[HEAD]** Year | City | Country | Nations **[ROW]** 1896 | Athens | …

**SQL Executor** — supervise → **Model**

# Experimental Result: **Effective Pre-training**

## WikiSQL (Weak)



72.4 — MAPO
74.8 — MeRL
83.9 — HardEM
83.6 — TAPAS
84.7 — GraPPa
89.6 — TAPEX

**+3.8**

## TabFact



65.1 — BERT
71.7 — Logic
72.3 — HeterTFV
74.4 — ProgVGAT
81.0 — TAPAS+
84.2 — TAPEX

**+3.0**

## SQA



44.7 — DynSP
55.1 — RA
67.2 — TAPAS
65.4 — SCoRE
71.0 — TAPAS+
74.5 — TAPEX

**+15.9**

## WikiTableQuestions



43.8 — MAPO
44.1 — MeRL
48.8 — TAPAS
52.3 — TaBERT
52.7 — GraPPa
57.5 — TAPEX

**+19.5**

# Experimental Result: **Efficient Pre-training**

### Fine-tuning Performance

| | | |
|---|---|---|
| TAPAS | TaBERT | TAPEX |
| 48.8 | 52.3 | 54.2 |

### Pre-training Corpus (Million)

| | | |
|---|---|---|
| TaBERT | TAPAS | TAPEX |
| 26.3 | 21.3 | 0.5 |

Compared with TaBERT, **2%** of corpus yields **2%** improvement!

# Experimental Analysis: **Larger is Better**

Scaling up the pre-training corpus generally brings positive effects.

# Experimental Analysis: **Fine-grained Analysis**

TAPEX significantly boosts the performance on all operators, implying that it does enhance BART's capabilities for joint reasoning over text and tables.

| Operator | Example Question | BART | TAPEX |
|---|---|---|---|
| **Select** | What is **the years won** for each team? | 41.3% | 64.8% (+23.5%) |
| **Filter** | How long did **Taiki Tsuchiya** last? | 40.1% | 65.7% (+25.6%) |
| **Aggregate** | What is the **amount of** matches drawn? | 26.9 % | 57.4% (+30.5%) |
| **Superlative** | What was the **last** Baekje Temple? | 46.3 % | 64.3% (+18.0%) |
| **Arithmetic** | What is the **difference** between White voters and Black voters in 1948? | 33.1 % | 53.5% (+20.4%) |
| **Comparative** | Besides Tiger Woods, what other player won **between 2007 and 2009**? | 30.0 % | 55.9% (+25.9%) |
| **Group** | What was score **for each** winning game? | 49.5 % | 66.7% (+17.2%) |

# Experimental Analysis: Complexity

Adding simpler SQL queries can improve performance on harder questions.

| Difficulty | Example SQL Query |
|---|---|
| Easy | SELECT Date<br>SELECT COUNT (Canal)<br>SELECT Name WHERE Age >= 28 |
| Medium | SELECT Region ORDER BY ID DESC LIMIT 1<br>SELECT COUNT (Tornadoes) WHERE Date = 1965<br>SELECT District WHERE District != "Tikamgarh" AND Agg = 0 |
| Hard | SELECT (SELECT COUNT( Distinct Area)) >= 5<br>SELECT COUNT (*) WHERE Result = "won" AND Year > 1987<br>SELECT Driver WHERE Manufacturer = "t-bird" ORDER BY Pos ASC LIMIT 1 |
| Extra Hard | SELECT COUNT (*) WHERE Position = 1 AND Notes = "110 m hurdles" AND Year > 2008<br>SELECT Nation WHERE Nation != "Japan" AND Gold = (SELECT Gold WHERE Nation = "Japan" )<br>SELECT Tournament WHERE Tournament IN ("oldsmar", "los angeles") GROUP BY Tournament ORDER BY COUNT (*) DESC LIMIT 1 |



Question Difficulty Level in Downstream

|  | BART | ≤ Easy | ≤ Medium | ≤ Hard | ≤ Extra Hard |
|---|---|---|---|---|---|
| Extra Hard | 27.5 | 28.3 | 32.5 | 40.8 | 42.5 |
| Hard | 40.0 | 42.6 | 53.1 | 58.8 | 60.2 |
| Medium | 34.4 | 38.2 | 56.2 | 57.3 | 56.9 |
| Easy | 57.4 | 63.9 | 70.2 | 70.2 | 71.7 |

SQL Difficulty Level in Pre-training

# Experimental Analysis: **Naturalness**

However, replacing SQL with NL does not benefit the pre-training, because the translated NL sentences contain noise.

| SQL Query | Translated NL Sentence | Faithfulness |
|---|---|---|
| **SELECT** Name **WHERE** Age **>=** 28 | Who is at least 28 years old? | ✓ |
| **SELECT MAX** (Pick#) | What was the last pick in the 1989 major league baseball draft? | ✗ |
| **SELECT** Driver **ORDER BY** Pos **DESC LIMIT 1** | What driver came in last place? | ✓ |
| **SELECT COUNT** (Competition) **WHERE** Notes **!=** 100 | How many competitions have no notes? | ✗ |
| **SELECT COUNT** (*) **WHERE** Result = "won" **AND** Year **>** 1987 | How many times did they win after 1987? | ✓ |
| **SELECT MAX** (Chart Position) **−** **MIN** (Chart Position) **WHERE** Release date **=** "july 21, 1995" | What is the difference between the chart position of july 21, 1995 and the chart position of july 22, 1995? | ✗ |
| **SELECT** Nation **WHERE** Nation **!=** "Japan" **AND** Gold **=** (**SELECT** Gold **WHERE** Nation **=** "Japan" ) | Which other countries had the same number of gold medals as Japan? | ✓ |
| **SELECT** Incumbent Electoral History **GROUP BY** Incumbent Electoral History **ORDER BY COUNT** (*) **DESC LIMIT 1** | Who has held the office the most? | ✗ |

# Take Away: Pre-training without Real Data

When performing continual pre-training, instead of mining a large noisy web corpus, we can also try to synthesize an accurate and small corpus.

# Take Away: Pre-training without Language Modeling

When performing continual pre-training, instead of performing the general-purpose language modeling, we can also try to simulate the specialized skill.

# POET: Reasoning Like Program Executor

*Xinyu Pi[1]   *Qian Liu[2]   Bei Chen[3]   Morteza Ziyadi[3]   Zeqi Lin[3]   Yan Gao[3]   Qiang Fu[3]

Jian-Guang Lou[3]  Weizhu Chen[3]

1 UNIVERSITY OF ILLINOIS URBANA-CHAMPAIGN   2 北京航空航天大学 BEIHANG UNIVERSITY   3 Microsoft

# Background: **Numerical Reasoning**

## Document

In **1517**, the seventeen-year-old King sailed to Castile. There, his Flemish court … In **May 1518,** Charles traveled to Barcelona in Aragon.

## Question

Where did Charles travel to first, Castile or Barcelona?

## Answer

Castile

# Method: SQL Execution Pre-training

Since SQL queries involve rich numerical operations, we hope it can be leveraged to enhance the numerical reasoning capability of models on documents.



Synthetic Pre-training Corpus

Realistic Downstream Datasets

# Method: SQL Execution for Different LMs

**Random Table**

| Year | City | Country |
|------|------|---------|
| 1896 | Athens | Greece |
| 1900 | Paris | France |
| … | … | … |
| 2008 | Beijing | China |

SELECT City … [HEAD] Year | City | Country

[ROW] 1896 | Athens  Result …

↑ Query Result Selection

**Encoder-Only LM**

**Model Input**

SELECT City … [HEAD]
Year | City | Country
[ROW] 1896 | Athens …

**Random SQL Query**

**SELECT** City**WHERE** Country **=** Greece
**ORDER BY** Year **ASC LIMIT 1**

**Encoder-Decoder LM**

↓ Query Result Generation

Athens

# Experimental Result: **Reasoning Transfer**

# Method: **Math Expression Calculation**

Observing the reasoning transfer from (SQL query, Database) to (Question, Passage), we propose a simplified method which leverages math expression for pre-training.



| | Pre-training | | Fine-tuning | |
|---|---|---|---|---|
| **Pre-trained LM for Textual Data** | → | **Pre-trained Reasoning LM** | → | **Fine-tuned Reasoning Model** |

**Input**: a + b - c [SEP] a = 2; d = 8; b = 5.2; c = 6.6; y = -12.5; x = 111; z = 999
**Output**: 3.4

Synthetic Pre-training Corpus

**Input**: Where did Charles travel to first, Castile or Barcelona? [Document]
**Output**: Castile

Realistic Downstream Datasets

F1 on DROP dataset based on BART

**69.2%** **78.1%**

# Experimental Analysis: **Performance Hurt on Other Tasks?**

**Small (<1%).** POET barely sacrifices the intrinsic understanding ability of language models.

# Experimental Analysis: **Benefit from Similarity of SQL to NL?**

**NO.** Randomly mapping SQL keywords to the "strange" tokens still works well.

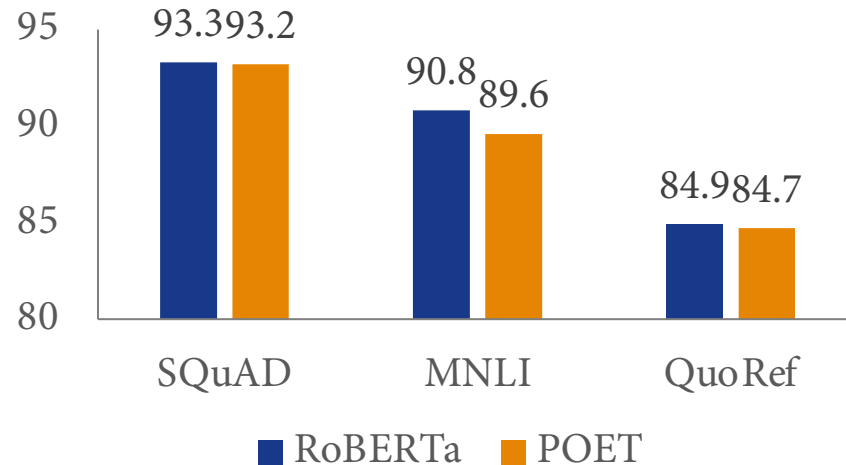SELECT City FROM table ORDER BY Population DESC LIMIT 1

*Fine-tune* ⟹ Which country contains the second largest part of the forest

**77.7%**

≋

unfocusedRange City guiIcon table externalToEVA awdownload Population ffffcc awdownloadclon 1

*Fine-tune* ⟹ Which country contains the second largest part of the forest

**76.9%**

# Experimental Analysis: **Pre-training on DROP Benefit SQL Execution?**

**Yes.** Pre-training on DROP leads to observably lower perplexity for SQL execution learning on both the train and dev sets.

# Experimental Analysis: **How Does it Work?**

**No answer.** But we can get some insights from the following analogy.

| Math Expression |
|:---:|
| x + y - z |

| Program |
|:---:|

≀≀

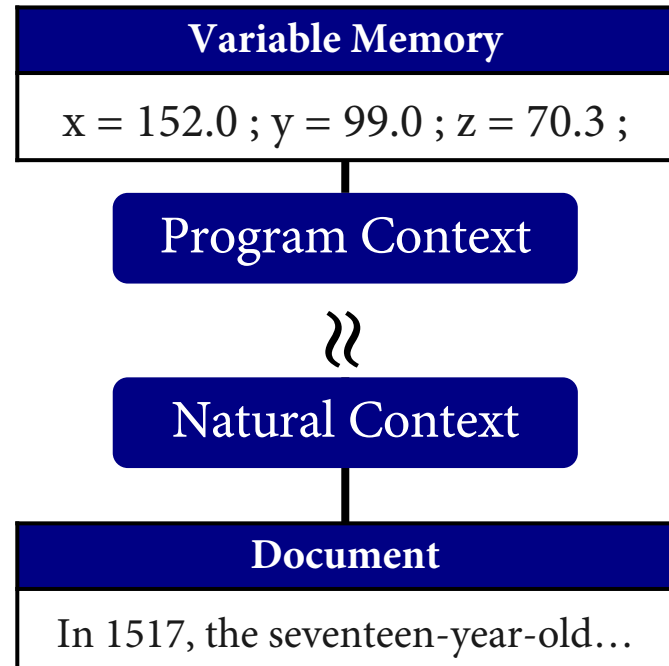| Natural Language |
|:---:|

| Question |
|:---:|
| Where did Charles travel to first? |

| Variable Memory |
|:---:|
| x = 152.0 ; y = 99.0 ; z = 70.3 ; |

| Program Context |
|:---:|

≀≀

| Natural Context |
|:---:|

| Document |
|:---:|
| In 1517, the seventeen-year-old… |

# Experimental Analysis: How Does it Work?

Without program context, the pre-training cannot work well.

| Math Expression |
|---|
| x + y - z |

Program

| Variable Memory |
|---|
| x = 152.0 ; y = 99.0 ; z = 70.3 ; |

Program Context

✔

| Math Expression |
|---|
| 2 + 8 − 6.6 |

Program

✘

# Take Away: Reasoning Transfer Occurs Across Modalities

Reasoning transfer occurs across modalities, and the analogy between pre-training and fine-tuning is important for the transference.

# LEMON: Language-Based Environment Manipulation via Execution-Guided Pre-training

Qi Shi[1]     Qian Liu[2]     Bei Chen[3]     Yu Zhang[1]     Ting Liu[1]     Jian-Guang Lou[3]

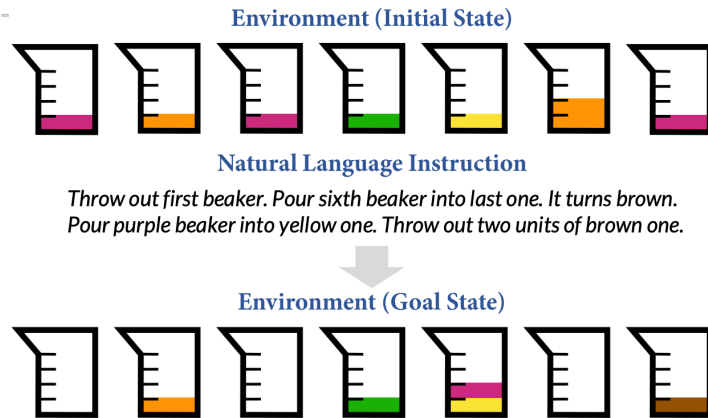[1] 哈尔滨工业大学 HARBIN INSTITUTE OF TECHNOLOGY     [2] 北京航空航天大学 BEIHANG UNIVERSITY     [3] Microsoft

# Background: **Language-Based Environment Manipulation**

Agents are required to manipulate the environments based on the natural language.
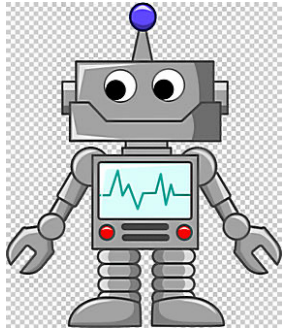
**Instruction Following**



**Environment (Initial State)**

**Natural Language Instruction**

*Throw out first beaker. Pour sixth beaker into last one. It turns brown.*
*Pour purple beaker into yellow one. Throw out two units of brown one.*

**Environment (Goal State)**

**Procedural Text Understanding**

| Paragraph (seq. of steps): | | Participants: | | | | |
|---|---|---|---|---|---|---|
| | | water | light | CO2 | mixture | sugar |
| **Roots absorb water from soil** | state0 | soil | sun | ? | - | - |
| **The water flows to the leaf.** | state1 | roots | sun | ? | - | - |
| **Light from the sun and CO2 enter the leaf.** | state2 | leaf | sun | ? | - | - |
| **The light, water, and CO2 combine into a mixture.** | state3 | leaf | leaf | leaf | - | - |
| **Mixture forms sugar.** | state4 | - | - | - | leaf | - |
| | state5 | - | - | - | - | leaf |

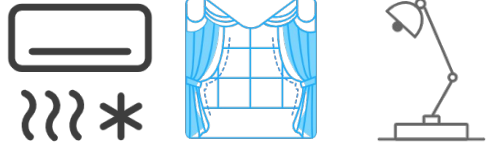**Time**

# Application: Language-Based Environment Manipulation



Swap the floor under the TV please.

Turn on the desk lamp, and turn off after 15 minutes.
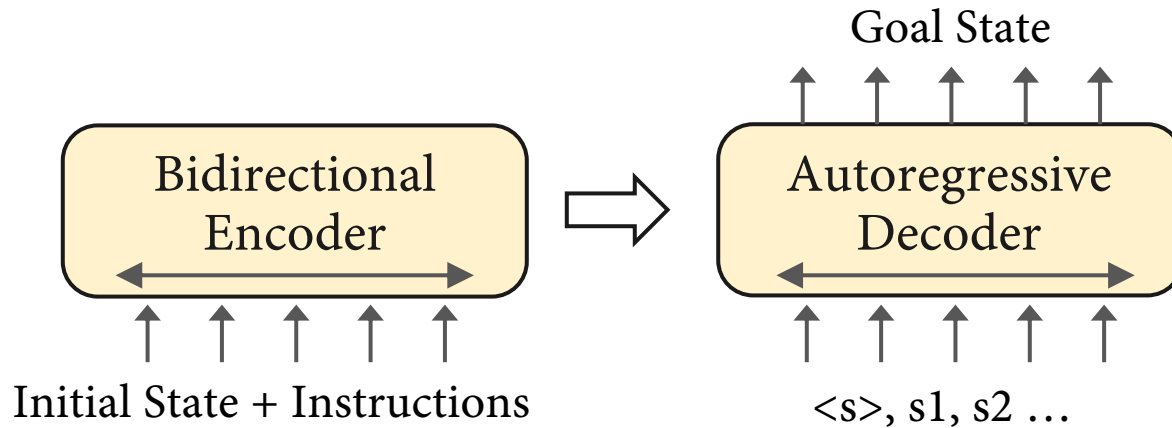
Jump on the box.

Agent Control

State Tracking

Virtual Interaction

# Preliminary: Generative Language Model Again

We formulate the task as a seq2seq paradigm, by leveraging generative PLMs (e.g., BART) to generate goal states directly.
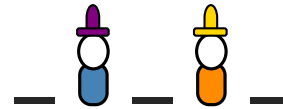
Goal State

Bidirectional Encoder ⟹ Autoregressive Decoder

Initial State + Instructions

<s>, s1, s2 …

# Challenge: **Spatial Reasoning**

Since pre-trained language models does not observe environments before, it is difficult for them to perform accurate spatial reasoning.
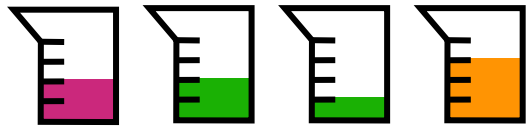
😳 What are these?

**water**  **light**  **carbon**

# Motivation: Environment Exploration by Actions

Synthesizing diverse actions to drive LMs familiar with environments.
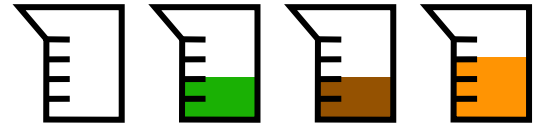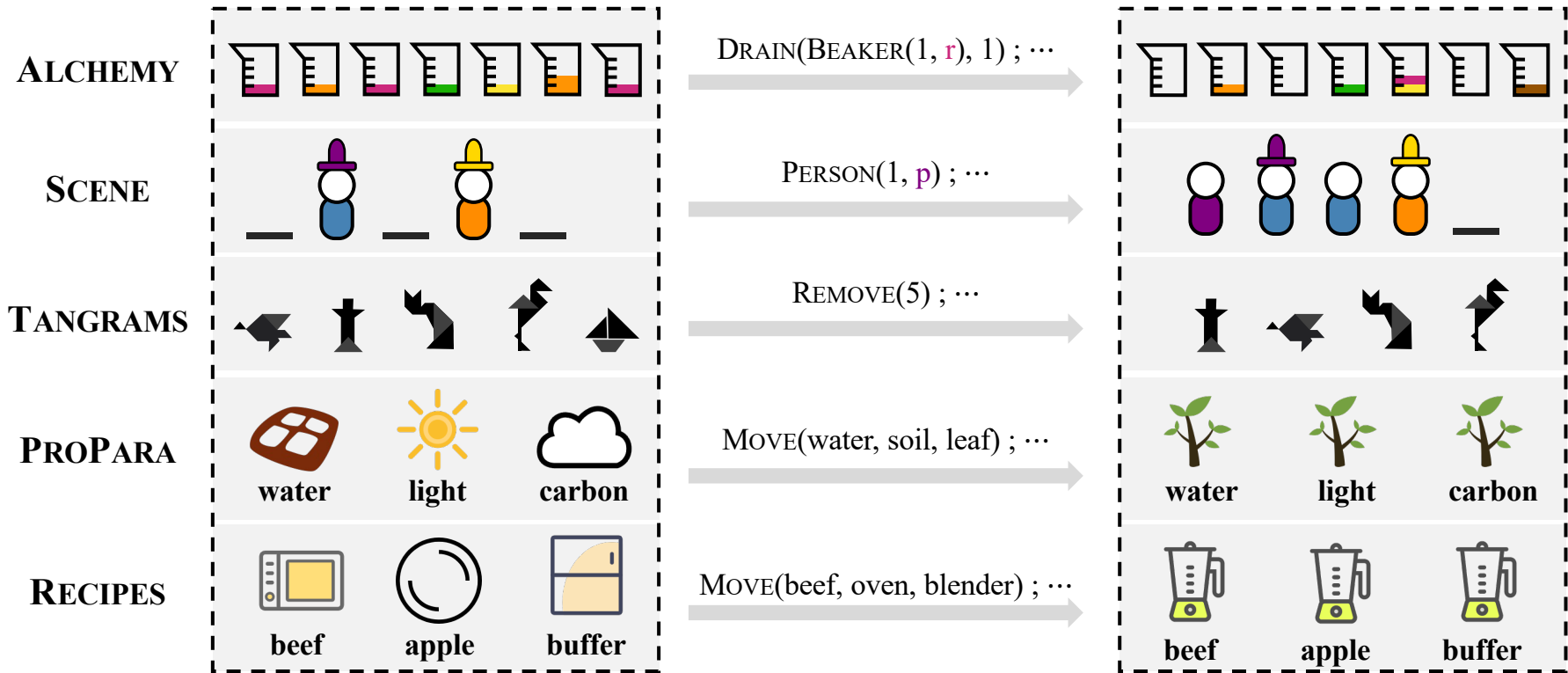


Environment (Initial State)

Action

Pour (Beaker (1), Beaker (2, g));
Drain (Beaker (3), $\frac{1}{3}$ );
Mix (Beaker (3));

LM

Environment (Goal State)

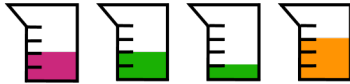# Method: Environment Exploration by Actions

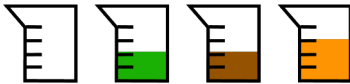# Method: **Environment Exploration by Actions**

## Pre-training

**Environment (Initial State)**

**Program**

POUR (BEAKER (1), BEAKER (2, g));
DRAIN (BEAKER (3), $\frac{1}{3}$ );
MIX (BEAKER (3));

**Environment (Goal State)**

## Fine-tuning
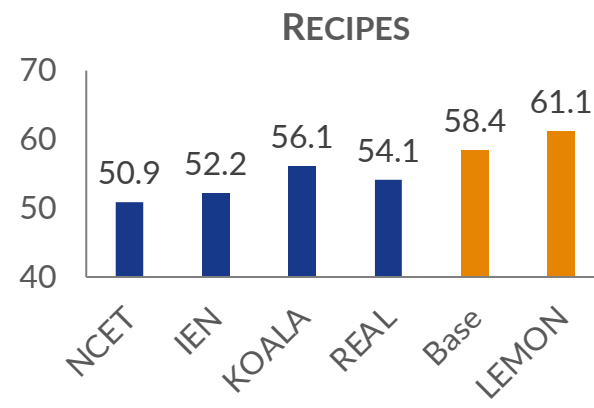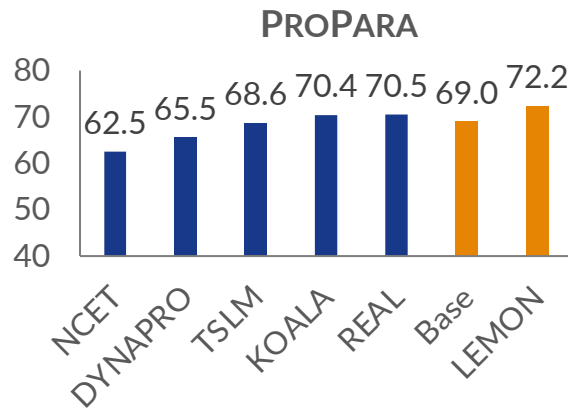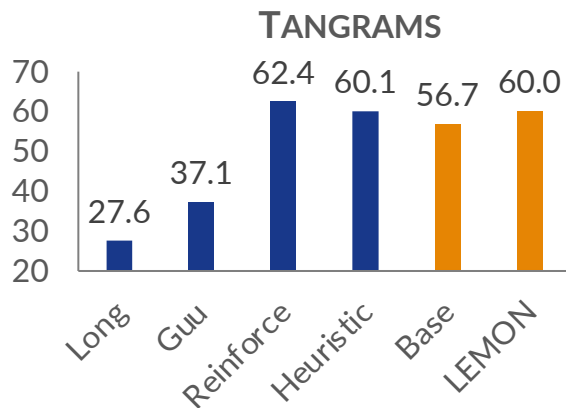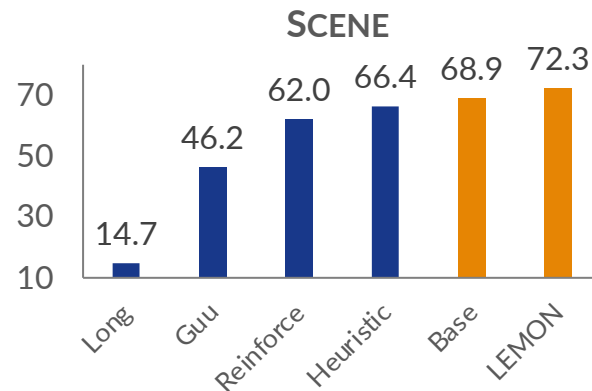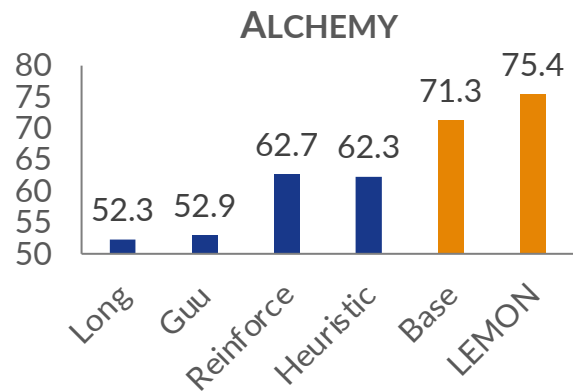
**Environment (Initial State)**

**Natural Language Instruction**

Throw out first beaker. Pour sixth beaker into last one. It turns brown. Pour purple beaker into yellow one. Throw out two units of brown one.
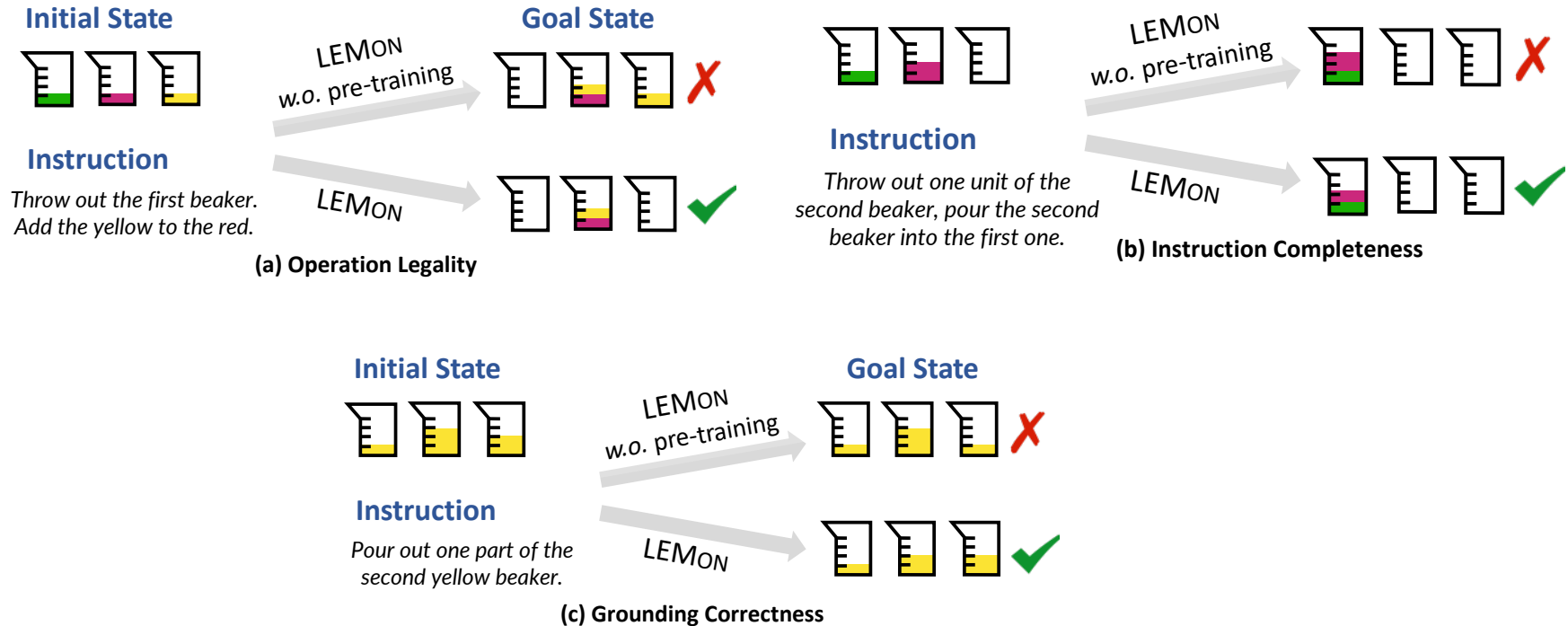
**Environment (Goal State)**
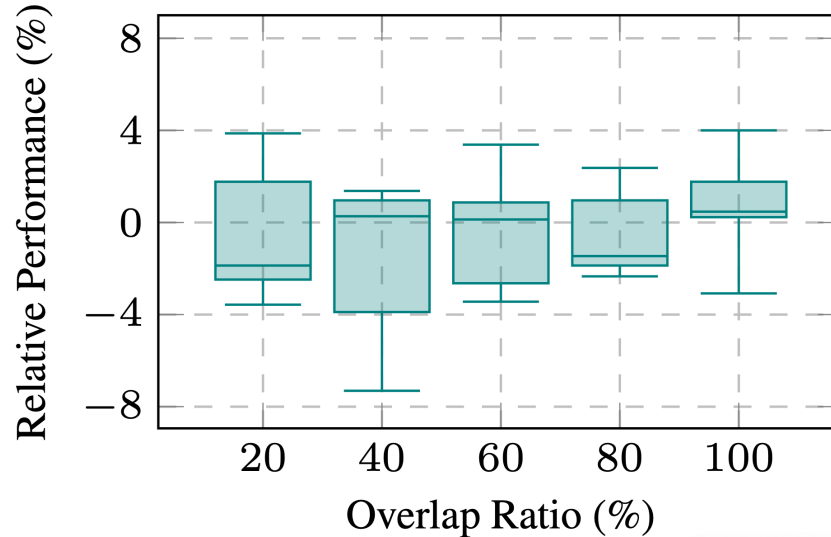
# Experimental Result: SOTA on Five Benchmarks



**ALCHEMY**

| Long | Guu | Reinforce | Heuristic | Base | LEMON |
|------|-----|-----------|-----------|------|-------|
| 52.3 | 52.9 | 62.7 | 62.3 | 71.3 | 75.4 |

**SCENE**

| Long | Guu | Reinforce | Heuristic | Base | LEMON |
|------|-----|-----------|-----------|------|-------|
| 14.7 | 46.2 | 62.0 | 66.4 | 68.9 | 72.3 |

**TANGRAMS**

| Long | Guu | Reinforce | Heuristic | Base | LEMON |
|------|-----|-----------|-----------|------|-------|
| 27.6 | 37.1 | 62.4 | 60.1 | 56.7 | 60.0 |

**PROPARA**

| NCET | DYNAPRO | TSLM | KOALA | REAL | Base | LEMON |
|------|---------|------|-------|------|------|-------|
| 62.5 | 65.5 | 68.6 | 70.4 | 70.5 | 69.0 | 72.2 |

**RECIPES**

| NCET | IEN | KOALA | REAL | Base | LEMON |
|------|-----|-------|------|------|-------|
| 50.9 | 52.2 | 56.1 | 54.1 | 58.4 | 61.1 |

# Experimental Analysis: What Does LEMON Learn?



(a) Operation Legality

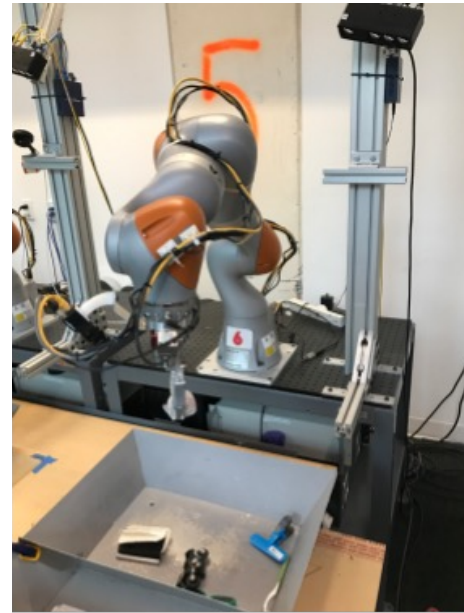(b) Instruction Completeness
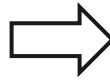
(c) Grounding Correctness

# Experimental Analysis: **Improvements from Leakage?**

**No**. The box plot of the relative performance (vertical axis) with respect to the overlap ratio (horizontal axis) indicates the independence.

# Take Away: Actions v.s. Simulation

Simulation to reality is a popular technique in autopilot. Actions can be regarded as kind of simulations which can facilitate the spatial reasoning in real space.

# Thanks & QA

**Qian Liu** (刘乾)

Research Scientist @ Sea AI Lab