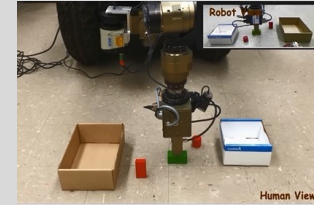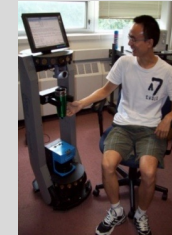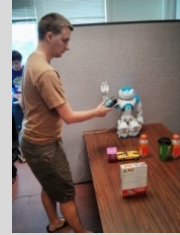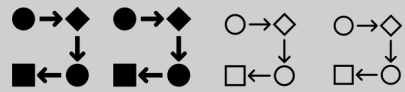# Large Pre-Trained Language Models for Physical Action Understanding and Planning

Shane Storks & Jianing "Jed" Yang

SLED Research Group @ University of Michigan

Oct. 21st, 2022

# Situated Language and Embodied Dialog
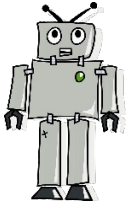


Mental Model
Background Knowledge

Language Communication

Mental Model
Background Knowledge

**Physical world**
- Objects, attributes, spatial relations
- Actions, states, goals
- Plans, Task structures

Planning, action

Perception, reasoning

Planning, action

Perception, reasoning

Language Grounding/Learning

Language Grounding/Learning

*(Slide from Joyce Chai)*

# Understanding Physical Causality

**6 -8 months**

Notice relationships between events. Perform basic actions to make things happen

**18 months**

Combine simple actions to make things happen. Change the way how they interact with the world to see how it changes the outcome

**36 months**

Make prediction about what may happen and reflect upon what caused something to happen

# Outline

1. Understanding the ability of large language models (LMs) to learn **verifiable physical commonsense reasoning**

2. Applying large LMs as a tool to inform **planning of physical actions**

# Motivation

- NLP tasks commonly boil natural language understanding (NLU) down to simple text classification tasks
    - Data bias and lack of transparency make it unclear whether underlying problems are truly solved
    - We want to examine system's underlying reasoning capability
- Tiered Reasoning for Intuitive Physics (`TRIP`) provides traces of a multi-tiered, human-annotated reasoning process:
    - Low-level, concrete physical states
    - High-level end task of plausibility classification

Storks, S., Gao, Q., & Chai, J. (2021). Tiered Reasoning for Intuitive Physics: Toward Verifiable Commonsense Language Understanding. Findings of EMNLP 2021.

# Tiered Reasoning for Intuitive Physics (`TRIP`)

## Story A

1. Ann sat in the chair.
2. Ann turned off the telephone.
3. Ann picked up a pencil.
4. Ann opened the book.
5. Ann wrote in the book.

## Story B

1. Ann sat in the chair.
2. Ann turned off the telephone.
3. Ann picked up a pencil.
4. Ann opened the book.
5. Ann heard the telephone ring.

**Which story is more plausible? A**

**Why not B?**

  **Conflicting sentences**: 2 → 5

  **Physical states:**

Powered(telephone) ⟶ ¬Powered(telephone)

Powered(telephone) ⟶ Powered(telephone)
Running(telephone)

Storks, S., Gao, Q., & Chai, J. (2021). Tiered Reasoning for Intuitive Physics: Toward Verifiable Commonsense Language Understanding. Findings of EMNLP 2021.

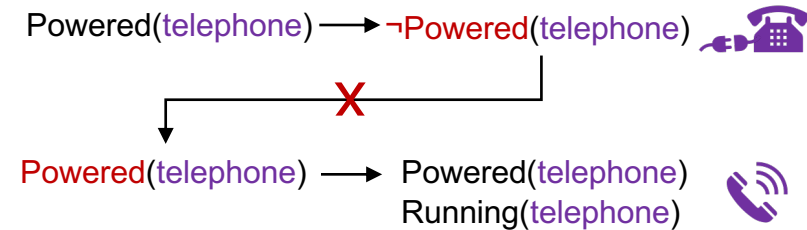# Evaluation Metrics

2. Ann turned off the telephone.

3. Ann picked up a pencil.

4. Ann opened the book.

5. Ann heard the telephone ring.

| Metric | Story Choice | Conflicting Sentences | Physical States |
|---|---|---|---|
| *Accuracy* | ✔ | | |
| *Consistency* | ✔ | ✔ | |
| *Verifiability* | ✔ | ✔ | ✔ |

Goal: Accuracy ≈ Consistency ≈ Verifiability

# Tiered Baseline



$$\mathcal{L} = \lambda_p \mathcal{L}_p + \lambda_f \mathcal{L}_f + \lambda_c \mathcal{L}_c + \lambda_s \mathcal{L}_s$$

Simply fine-tuning a pre-trained LM on the end task (plausibility prediction) **can achieve up to 97% accuracy**.

# However…

# Results of T-NLR (Large) on TRIP



Chart showing percentages (%) for Accuracy, Consistency, and Verifiability across four conditions:

| | Accuracy | Consistency | Verifiability |
|---|---|---|---|
| All Losses | 69.6 | 2.8 | 0.3 |
| Omit Story Choice Loss | 75.8 | 7.1 | 2.8 |
| Omit Conflict Detection Loss | 66.8 | 1.2 | 0.6 |
| Omit State Classification Losses | 74.5 | 2.2 | 0.0 |

Legend: ■ Accuracy ■ Consistency ■ Verifiability

Gao, J & Tiwary, S. (2021). Efficiently and effectively scaling up language model pretraining for best language representation model on GLUE and SuperGLUE.

# Error Distribution of T-NLR



Inaccurate

Accurate

Consistent but not verifiable!

Correct and entirely verifiable!

Correct states, but unsuccessful conflict detection. 🤔

Correct, but entirely unverifiable!

24.2%

4.3% 2.8%

6.2%

22.4%

1.2%

0.6%

62.4%

75.8%

SC: √  PS: √
SC: √  PS: X
SC: X  PS: √
SC: X  PS: X

SC: sentence conflict
PS: physical states

# Embodied Task Reasoning



Household Task Domain

Instructions

Task Learning

Navigation

Manipulation

Compositional

Generalize to unseen scenes

- Pick and place
- Do cleaning
- Prepare coffee
- Cook food ...

# Task Reasoning in Simulated Environment

**ALFRED** (**A**ction **L**earning **F**rom **R**ealistic **E**nvironments and **D**irectives)



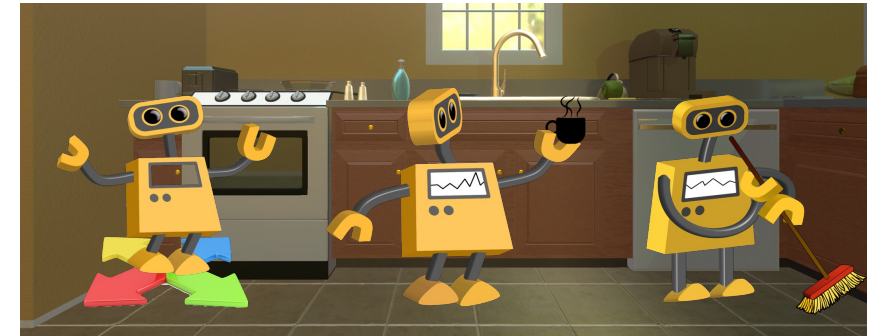ALFRED

25k human language annotations

8k demonstrations

2k unique tasks

120 scenes

56 object classes

26 receptacle classes

High-level Goal Directive

Low-level Instructions

Visual Navigation

Object Interaction

Shridhar et al. ALFRED: A Benchmark for Interpreting Grounded Instruction for Everyday Tasks, CVPR 2020

# Hierarchical Task Learning with Unified Transformers (HiTUT)

Previous Work



(Shridhar et al., 2020; Pratap Singh et al., 2020; Storks et al., 2021)

# Hierarchical Task Learning with Unified Transformers (HiTUT)



Previous Work

**End-to-End Modeling**

Goal Directive | Sub-Goal Instructions

(Shridhar et al., 2020; Pratap Singh et al., 2020; Storks et al., 2021)

HiTUT

Sub-Goal Instructions

Goal Directive

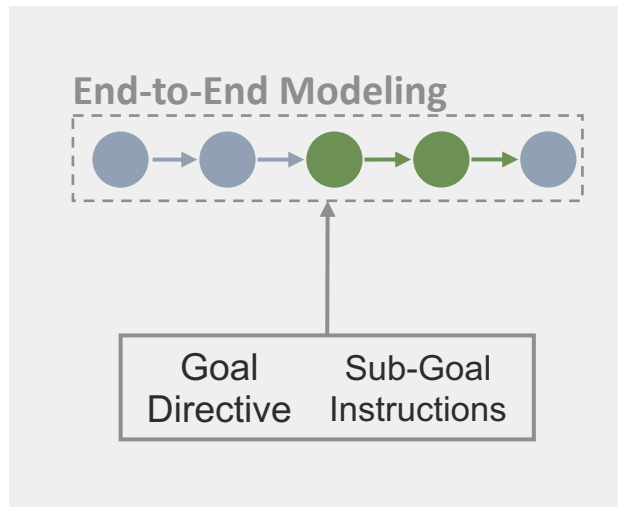Primitive Actions

Sub-Goals

Navigation Sub-goals/Actions

Manipulation Sub-goals/Actions

Unified Transformers (e.g. BERT)

Mask Selection   Mask   Type   Arg

Sum over heads & Softmax   Object Detector   FC & Softmax

K  K  K   **BERT**   Q

LN   FC & LN   FC & LN   FC & LN

Emb   FC   Word Embedding

position   class

Object Detector   Tokenize   Predicate to Word

Posture Feature (Navi. only)   Visual Observation $\mathcal{V}$   Language Instruction $\mathcal{L}$   Predicate History $\mathcal{P}$

16

# Hierarchical Task Learning with Unified Transformers (HiTUT)

**Goal Directive** *Place a cleaned mug in the coffee machine.*

# Results: Better Generalization in Unseen Environment



**Task Goal:**
Put two books on the desk.

SG1: Goto(Book)
A1: LookDown

Diagnosis Results

Example of how backtracking helps the agent recover from execution errors.

# DANLI: A Deliberative Neuro-symbolic Agent

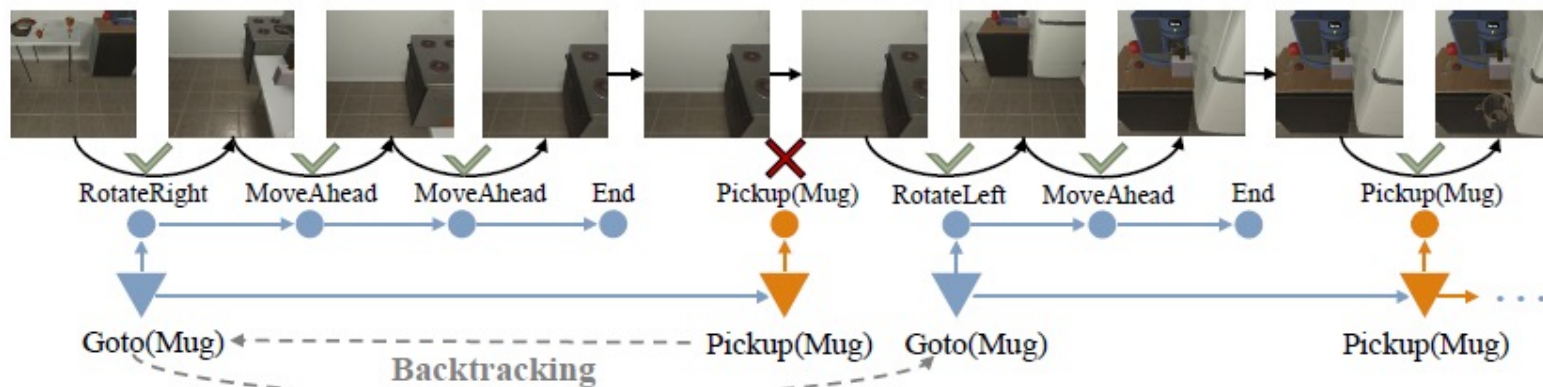- Leverages neural models to predict subgoals from dialog history, and constructs a 3D voxel map representation from agent's ego-centric vision

- Support symbolic reasoning and proactive planning using a PDDL-based online planning algorithm

DANLI: Deliberative Agent for Following Natural Language Instructions. EMNLP 2022.

# LLM for Long-horizon Planning: Pipeline



Goal Instruction

*"Make a cup of coffee"*

Retrieve

Plan Bank

```
---------- GIVEN EXAMPLE ----------
Task: Make toast
Step 1: Walk to dining room
Step 2: Walk to freezer
Step 3: Find freezer
Step 4: Open freezer
Step 5: Find food bread
Step 6: Grab food bread
Step 7: Close freezer
Step 8: Find toaster
Step 9: Plug in toaster
Step 10: Put food bread on toaster
Step 11: Switch on toaster
```

Prompt

LLM

Environment Feedback

Actuate

Raw Action
*"grab cup"*

Translated Action

*Pickup(Cup)*

SentenceBERT

# Evaluation: Zero-Shot Action Prediction

| | GPT-2 (1.5B) | Turing-7B |
|---|---|---|
| Action Error Rate (↓) | 89.4% | **40.6%** |
| Acc. (Type+Arg) (↑) | 1.5% | **3%** |
| Acc. (Type) (↑) | **25.2%** | 23.3% |
| Acc. (Arg) (↑) | 4.5% | **8.7%** |

# Qualitative Examples: GPT-2 & Turing-7B

*Prompt: "<Example_Tasks> Task: Do the laundry, Step 1: _____"*

| GPT-2 (1.5B) | | | Turing (7B) | | |
|---|---|---|---|---|---|
| | **GPT-2 action** | **ALFRED action** | | **Turing action** | **ALFRED action** |
| **0** | walk to kitchen | Open drawer | **0** | walk to bedroom | Go to bed |
| **1** | pull out clothes | Move dresser | **1** | wake up | Pick up alarm clock |
| **2** | move closet | Move shelving unit | **2** | plug in alarm clock | Put alarm clock on box |
| cut off action | move bookshelf | Move shelving unit | **3** | go to bathroom | Go to toilet |
| | | | **4** | pick up toilet brush | Pick up toilet paper |
| | | | **5** | put toilet paper on box | Put toilet paper on drawer |
| | | | **6** | pick up towel | Pick up towel |
| | | | **7** | put towel on rack | Put towel on towel holder |
| | | | **8** | pick up washcloth | Pick up cloth |
| | | | **9** | put cloth on clothesline | Put cloth on tv stand |
| | | | **10** | pick up clothes | Pick up cloth |
| | | | **11** | put cloth on coffee table | Put cloth on coffee table |
| | | | **12** | pick up cloth | Pick up cloth |
| | | | **13** | put cloth on kitchen counter | Put cloth on counter top |
| | | | **14** | pick up cloth | Pick up cloth |
| | | | **15** | put cloth on floor | Put cloth on desk |
| | | | **16** | pick up cloth | Pick up cloth |
| | | | **17** | put cloth on floor | Put cloth on desk |

| | |
|---|---|
| repetitive | |
| mistranslation | |

24

# Qualitative Examples: Turing-7B & Turing-530B

*Prompt: "<Example_Tasks> Task: Make breakfast, Step 1: _____"*

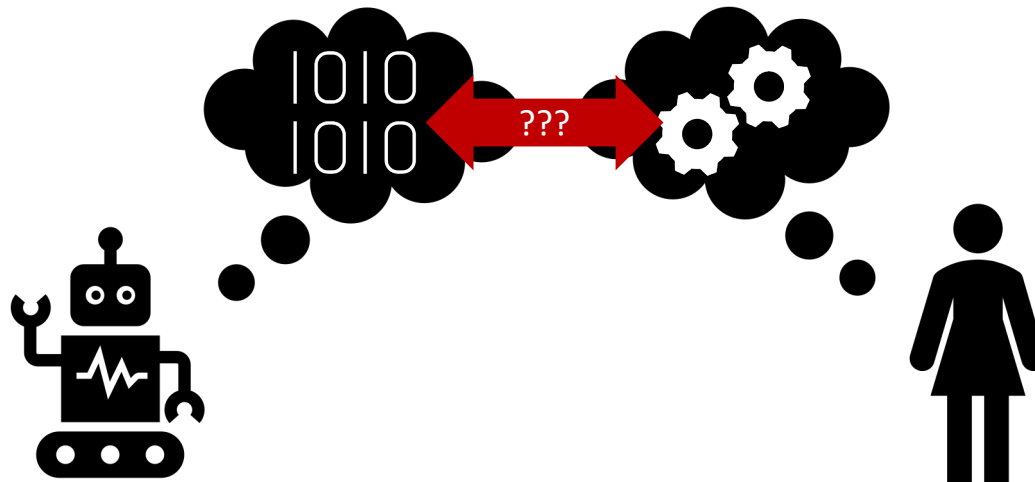| Turing (7B) | | | Turing (530B) | | |
|---|---|---|---|---|---|
| | **Turing action** | **ALFRED action** | | **Turing action** | **ALFRED action** |
| 0 | walk to dining room | Go to dining table | 0 | make toast | Toggle toaster |
| cut off action | put food on table | put plate on dining table | 1 | walk to dining room | Go to dining table |
| | | | 2 | find knife | Go to knife |
| | | | 3 | pick up knife | Pick up knife |
| | | | 4 | go to food bread | Go to bread |
| | | | 5 | cut bread | Slice bread |
| | | | 6 | place bread on plate | Put bread on plate |
| | | | 7 | go to cupboard | Put cup on cabinet |
| | | | 8 | pick up cup | Pick up cup |
| | | | 9 | fill cup with water | Fill watering can |
| | | | cut off action | water plants | Fill watering can |

| | |
|---|---|
| (yellow) | repetitive |
| (pink) | mistranslation |

# Summary

- While large LMs (such as T-NLR) make some steps toward coherent reasoning for NLU, more work is needed toward neuro-symbolic reasoning pipelines for teaching systems how to reason about the physical world.

- Large generative LMs (such as T-NLG) demonstrates some initial capability of zero-shot task planning, but still has large gap compared to fine-tuned LMs. More work is needed for translating and grounding LLM outputs to unseen task domains.

Coalescing Global and Local Information for Procedural Text Understanding. COLING 2022.
Language Models as Zero-Shot Planners: Extracting Actionable Knowledge for Embodied Agents. 2022.
Do As I Can, Not As I Say: Grounding Language in Robotic Affordances. 2022.
Inner Monologue: Embodied Reasoning through Planning with Language Models. 2022.

# Future Work

- Commonsense reasoning with large generative LMs
  - Analogy and relational reasoning
  - Generalized physical commonsense reasoning
- Action planning with large generative LMs
  - Close-loop planning utilizing environmental and interactive feedbacks

Shane Storks

@shanestorks

Jianing "Jed" Yang

@jed_yang

Qianqi Yan

Wenfei Tang

Joyce Y. Chai

33