# Towards a Robotics Foundation Model

Keerthana Gopalakrishnan March 3, 2023





# Agenda

- **01** The Ingredients for a Robotics Foundation Model
- 02 RT-1
- 03 SayCan
- 04 Inner Monologue
- 05 NL Map Saycan
- **06** MOO
- 07 What's next?

#### Why a Foundation Model for Robotics?

- Foundation models enable emergent capabilities and homogenization
  - *emergent capabilities:* emergence of more complex behavior not present in smaller models
  - *homogenization*: generalization to combinatorially many downstream use cases
  - More is Different for Al, Emergence in LLMs



#### Why a Foundation Model for Robotics?

- Foundation models enable emergent capabilities and homogenization
  - *emergent capabilities:* emergence of more complex behavior not present in smaller models
  - *homogenization*: generalization to combinatorially many downstream use cases
  - More is Different for AI, Emergence in LLMs
- An "emergent capabilities" curve might be *required* for robotics to be useful



#### Why not a Foundation Model for Robotics?



# What are the ingredients for a Robotics Foundation Model?



#### Ingredient #1: Design Principles of ML Scaling

- High-capacity architectures, ie. self-attention
- Scaling params and compute and corpus size (tokens)
- Dataset size matters more than quality





Andrej Karpa... @ @karpa... • Oct 19, 2022 •••
The Transformer is a magnificient neural network architecture because it is a general-purpose differentiable computer. It is simultaneously:
1) expressive (in the forward pass)
2) optimizable (via backpropagation+gradient descent)

3) efficient (high parallelism compute graph)

#### Ingredient #2: Proliferation of Internet-scale Models

- Generative models in {language, coding, vision, audio, ...} experience emergent capabilities
- Proliferation + acceleration means these models will get better "on their own" over time



Sources: "Compute trends across three eras of machine learning", by J. Sevilla et al., arXiv, 2022; Our World in Data

#### Ingredient #3: Robotics moves from Online to Offline

• Foundation models train on huge amounts of diverse offline datasets



The Pile dataset (Eleuther)

Dataset	# English Img-Txt Pairs			
Public Datasets				
MS-COCO	330K			
CC3M	3M			
Visual Genome	$5.4\mathrm{M}$			
WIT	$5.5\mathrm{M}$			
CC12M	12M			
RedCaps	12M			
YFCC100M	$100 \mathrm{M}^2$			
LAION-5B (Ours)	2.3B			
Private Datasets				
CLIP WIT (OpenAI)	400M			
ALIGN	1.8B			
BASIC	6.6B			

LAION-5B (Laion)

### Ingredient #3: Robotics moves from Online to **Offline**

Foundation models train on huge amounts of diverse offline datasets 



The Pile dataset (Eleuther)

Dataset # English Img-Txt Pairs Public Datasets MS-COCO 330K CC3M 3MVisual Genome 5.4MWIT 5.5MCC12M 12M RedCaps 12M  $100 {\rm M}^2$ YFCC100M LAION-5B (Ours) 2.3B**Private Datasets** CLIP WIT (OpenAI) 400M ALIGN 1.8IBASIC 6.6I LAION-5B (Laion) DETOUR Robot learning classically does a lot of online on-policy learning (RL!) but winds have been shifting



"How do we do end-to-end robot learning in the real world?"





learning in the real world?"

**#1**: Some methods plateaued at 50-70% success rates **#2**: Some methods required very specific data distributions

How do we satisfy both: #1 Solve many tasks (>90%) #2 Improve with "scalable" data

"How do we scale to many tasks in more complex scenarios?"

Spring 2022







We went from focusing on online methods to focusing on offline methods. We decoupled data generation from data consumption.

**Ingredients** 

<u>Lessons</u>











combine <u>large diverse offline datasets</u> with <u>high-capacity architectures</u> by using <u>language as a universal glue</u>

# Agenda

**01** The Ingredients for a Robotics Foundation Model

02 RT-1

- 03 SayCan
- 04 Inner Monologue
- 05 NL Map Saycan
- **06** MOO
- 07 What's next?

### Scaling Multi-task Imitation from First Principles



#### Scaling Multi-task Imitation from First Principles





- Tokenized input and outputs
- Decoder only transformer, sparse categorical entropy objective
- Image tokenizer: Pre-trained film efficient net backbone
- Token learner for compression/ faster inference



### RT-1 Design Takeaways

- Inference budget: 100ms (3Hz control)
- 6-image history
- TokenLearner subsamples 81 image patches into 8 image patches
- Action discretization to 256 bins
- Model Size: 35M parameters

### Pushing the limit: adding visual and semantic diversity

- RT-1 has the best performance compared to baselines in seen/unseen
- RT-1 is robust to variations in backgrounds, distractors



#### Pushing the limit: change multiple factors of variation





### Pushing the limit: Training on Diverse Data Distributions



### Scaling with data

- Reducing data size reduces performance and generalization
- Task diversity >> data size



% of Data

# Agenda

- **01** The Ingredients for a Robotics Foundation Model
- 02 RT-1
- 03 SayCan
- 04 Inner Monologue
- 05 NL Map Saycan
- **06** MOO
- 07 What's next



2016 - 2021

"How do we do end-to-end robot learning in the real world?"

#### 2021 - 2022

"How do we scale to many tasks with demonstrations?"

#### 2022 - 2023

"How can we leverage foundation models to accelerate robotics?"

#### Language models for robotics

Problem 1: Our robots can only do a fixed number of commands.

We need to get LLMs to speak "robot language"!

Problem 2: LLMs aren't grounded in the real-world. They can't "see".

We need to ground LLMs in robotic affordances!



#### SayCan: LLMs for robotics and robotics for LLMs



Possible Tasks with Affordances



#### say-can.github.io

### **Experiment Overview**

84% planning rate

- 74% execution rate
- 101 long-horizon instructions

10+ navigation and manipulation skills in a row Without grounding nearly halves performance

Instruction Family	Num	Plan	Execution
Natural Language Single Primitive	15	100%	100%
Natural Language Nouns	15	67%	47%
Natural Language Verbs	15	100%	93%
Structured Language	15	93%	87%
Embodiment	11	64%	55%
Crowd Sourced	15	87%	87%
Long-Horizon	15	73%	47%
Total	101	84%	74%



### Baly Maßay Can

Planning Performance

### Robotics performance scales with better LLMs!

FLAN-SayCan PALM-SayCan 100% 75% 50% 25% 0% 10 50 100 500 Model Size (B)



PaLM-SayCan

### Robotics performance scales with better LLMs!

- Chain-of-thought prompting
- Better prompts

#### Solves all kinds of queries:

- I'm thirsty
- Bring me two different sodas
- I left out a coke, apple, and water, can you throw them away and then bring me a sponge to wipe the table?
- 你能给我拿杯可乐吗?

Human Input: How would you bring me a fruit that is not an apple? Model Output: **Explanation:** The user has requested me to bring a fruit that is not an apple. The banana is a fruit that is not an apple, I will bring the user a banana. Robot: I would 1. find a banana 2. pick up the banana 3. go to the table 4. put down the banana 5. done.

# Agenda

- **01** The Ingredients for a Robotics Foundation Model
- 02 RT-1
- 03 SayCan
- 04 Inner Monologue
- 05 NL Map Saycan
- **06** MOO
- 07 What's next

#### Language as connective tissue for robotics



innermonologue.github.io

#### Inner Monologue

- Close the loop further with the scene
- Scene understanding

Kitchen Mobile Manipulation (Real)

• And human feedback for free

Human: Bring me a drink from the table Robot: go to the table Scene: I see a coke and lime soda Robot: What kind of drink would you like? Human: Something with caffeine Robot: Pick up coke Success Detector: False Robot: Pick up coke Robot: Bring it to you Robot: Put down the coke Robot: Done



#### **Object Recognition**

**Success Detection** 

# Agenda

- **01** The Ingredients for a Robotics Foundation Model
- 02 RT-1
- 03 SayCan
- 04 Inner Monologue
- 05 NL Map Saycan
- **06** MOO
- 07 What's next

#### Language as connective tissue for robotics





Queryable, natural language map provides world context to the planner.



nlmap-saycan.github.io







Use Foundation Models for Importing Common Sense

"The bottleneck for robotics is high level semantic planning"

#### - Robotics Professor in 2018

"No." - LLMs in 2022

Lesson #2: Leverage Foundation Models If language is the universal API of the system, then foundation models can inject common sense

# Agenda

- **01** The Ingredients for a Robotics Foundation Model
- 02 RT-1
- 03 SayCan
- 04 Inner Monologue
- 05 NL Map Saycan
- 06 MOO
- 07 What's next?

#### Two Axis for Generalization

Tasks

Examples:

Pick up X

Move X near Y

Open the lid on X

Place X into the Y

Objects

Examples:

Pick up the apple core

Place the screws into the box

Move the coffee near the human

Open the lid of the large box.

# Vision Language Models can zero-shot detect almost any object.



(Real detections from OWL-ViT)

#### Generalizing to any object type



#### Interface is natural language

Extract objects from phrase and feed them to OWL-ViT, feed detections as input along with "task embedding."



Mask

#### Train on enough objects the policy learns generalizable skills



#### **Distribution of Objects**

We added small amounts of "pick" skill data for 100 additional, diverse objects.



#### Main Results

	Pick		Other skills	
Method	Seen objects	Unseen objects	Seen objects	Unseen objects
RT-1 (our data) [2]	54	25	50	50
RT-1 (original data)	$31^{1}$	38	$17^1$	13
VIMA-like [13]	62	50	50	25
MOO (ours)	92	75	83	75

#### **Bonus: Better Environment Generalization**



Method	Open-World Objects	Challenging Textures	New Environments
RT-1 (our data) [2]	17	7	29
VIMA-like [13]	50	7	7
MOO (ours)	67	50	43

TABLE II: Robustness evaluations for novel use cases. MOO is able to handle new objects, textures, and environments with substantially greater success than prior methods.

#### Limitations

- Still struggles with objects far outside training distribution.
- Failures often do plausible things and exhibit retry behavior
- Learning new action primitives is still very hard.



#### Takeaways

Lesson #2: Leverage Foundation Models

**Lesson #3:** Offline Robot Learning • Use VLMs for generalization

- Import internet-scale priors into offline datasets
- When datasets are large enough, some label noise is okay

# Agenda

- 01 The Ingredients for a Robotics Foundation Model
- 02 RT-1
- 03 SayCan
- 04 Inner Monologue
- 05 DIAL
- 06 What's Next?



<u>Recipe</u>

combine <u>large diverse offline datasets</u> with <u>high-capacity architectures</u> by using <u>language as a universal glue</u>

combine <u>large diverse offline datasets</u> with <u>high-capacity architectures</u> by using <u>language as a universal glue</u>



Component	Method
Skill Learning	<u>RT-1</u>
Planning	<u>SayCan, Inner Monologue</u>
Low-level Control	Code as Policies
Data Augmentation	DIAL, ROSIE
Object-centric Representations	NLMap, MOO



2016 - 2021

"How do we do end-to-end robot learning in the real world?"

#### 2021 - 2022

"How do we scale to many tasks with demonstrations?"

#### 2022 - 2023

"How can we leverage foundation models to accelerate robotics?"

#### **Open Research Directions**

- Bottleneck is still on skill learning
- Leveraging Bitter Lesson 2.0 is non-trivial
- How do we collect diverse and useful data more scalably?
- What algorithms are absorbent "data sponges"?
- Transfer from human embodiment



### Thank you!

<u>say-can.github.io</u> <u>innermonologue.github.io</u> robotics-transformer.github.io <u>robot-moo.github.io</u> <u>nlmap-saycan.github.io</u>

