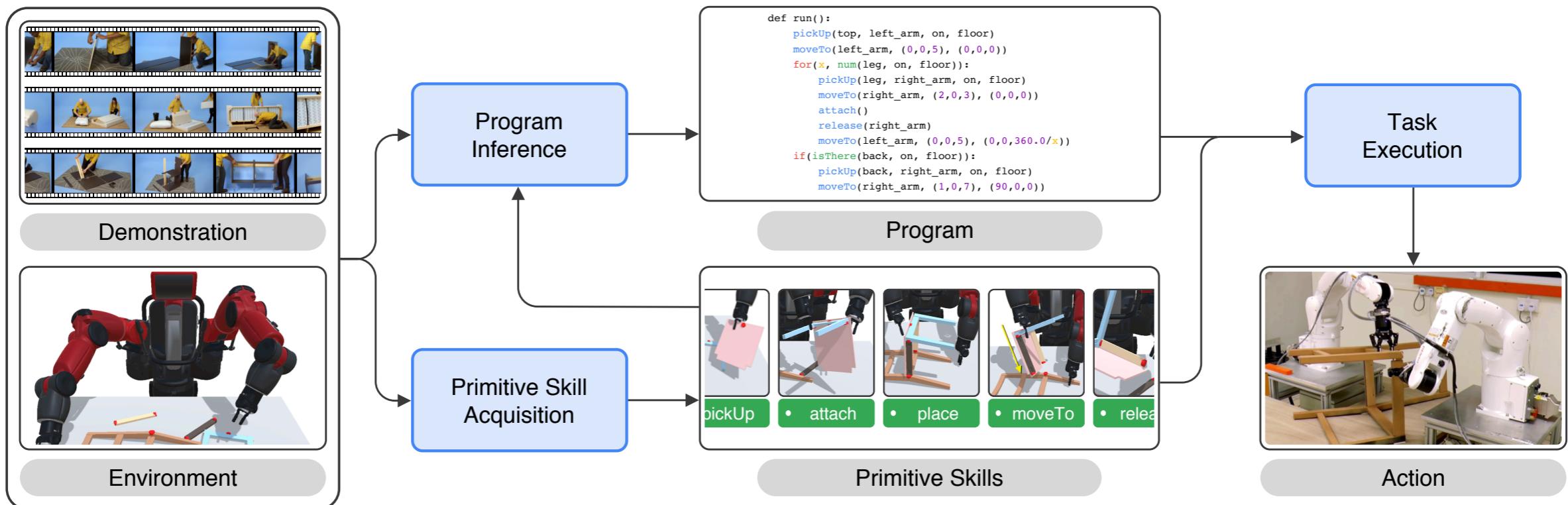
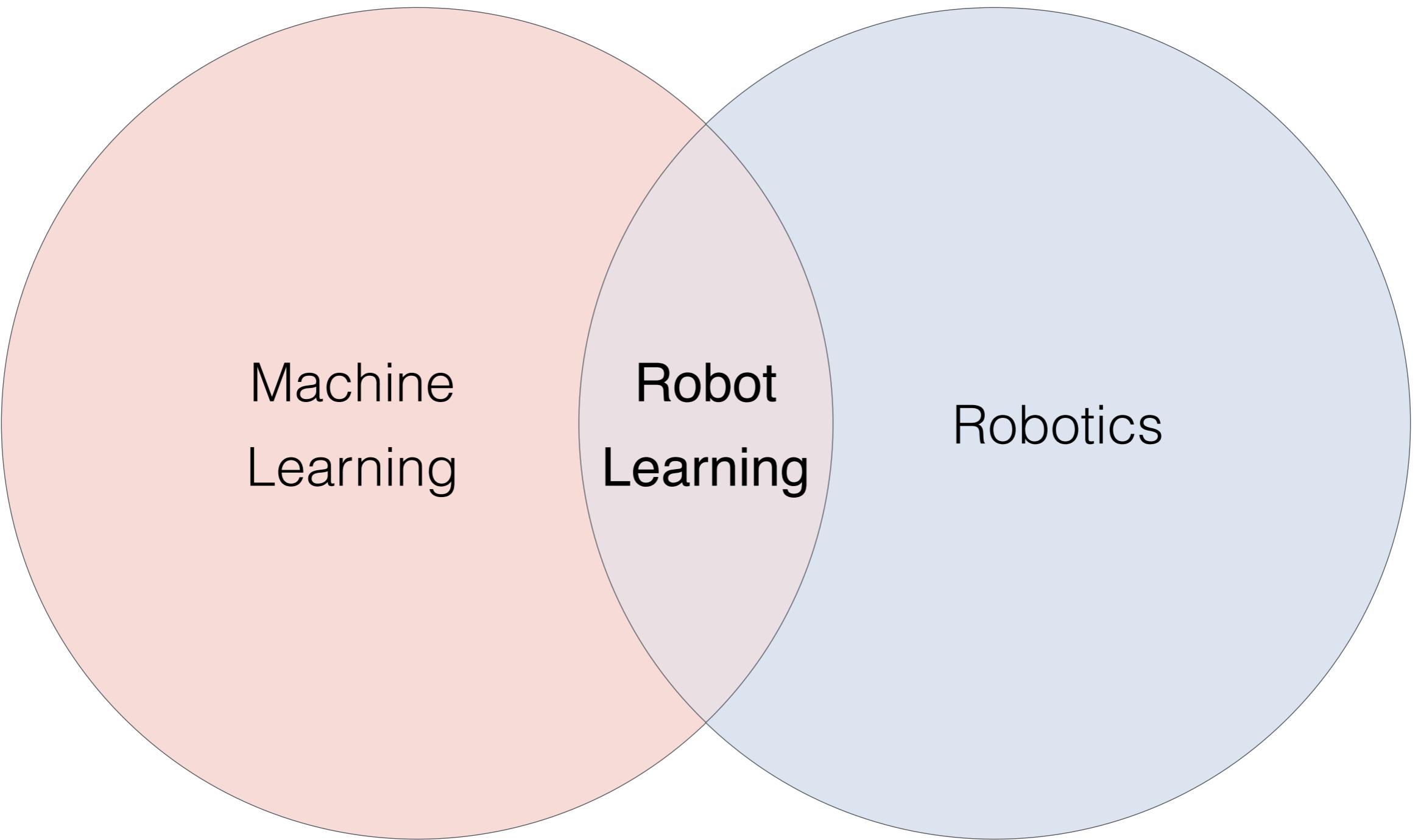


Program-Guided Framework for Interpreting and Acquiring Complex Skills with Learning Robots

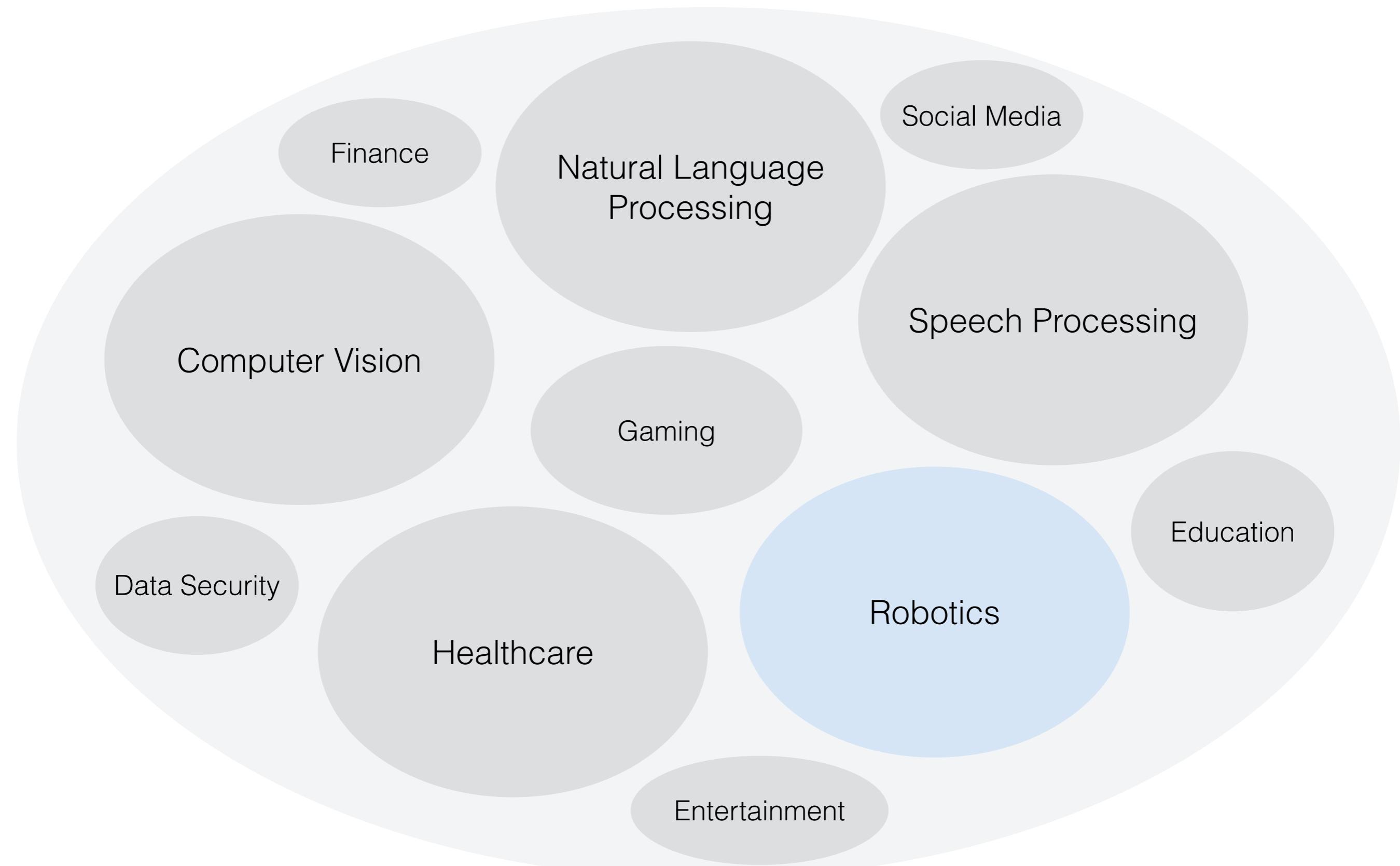


Shao-Hua Sun (孫紹華)

Ph.D. candidate in Computer Science
at the University of Southern California (USC)



Applications of AI



Learning Problems

Supervised Learning

Active Learning

Meta-Learning

Unsupervised Learning

Reinforcement Learning

Meta-Reinforcement Learning

Sparse Dictionary Learning

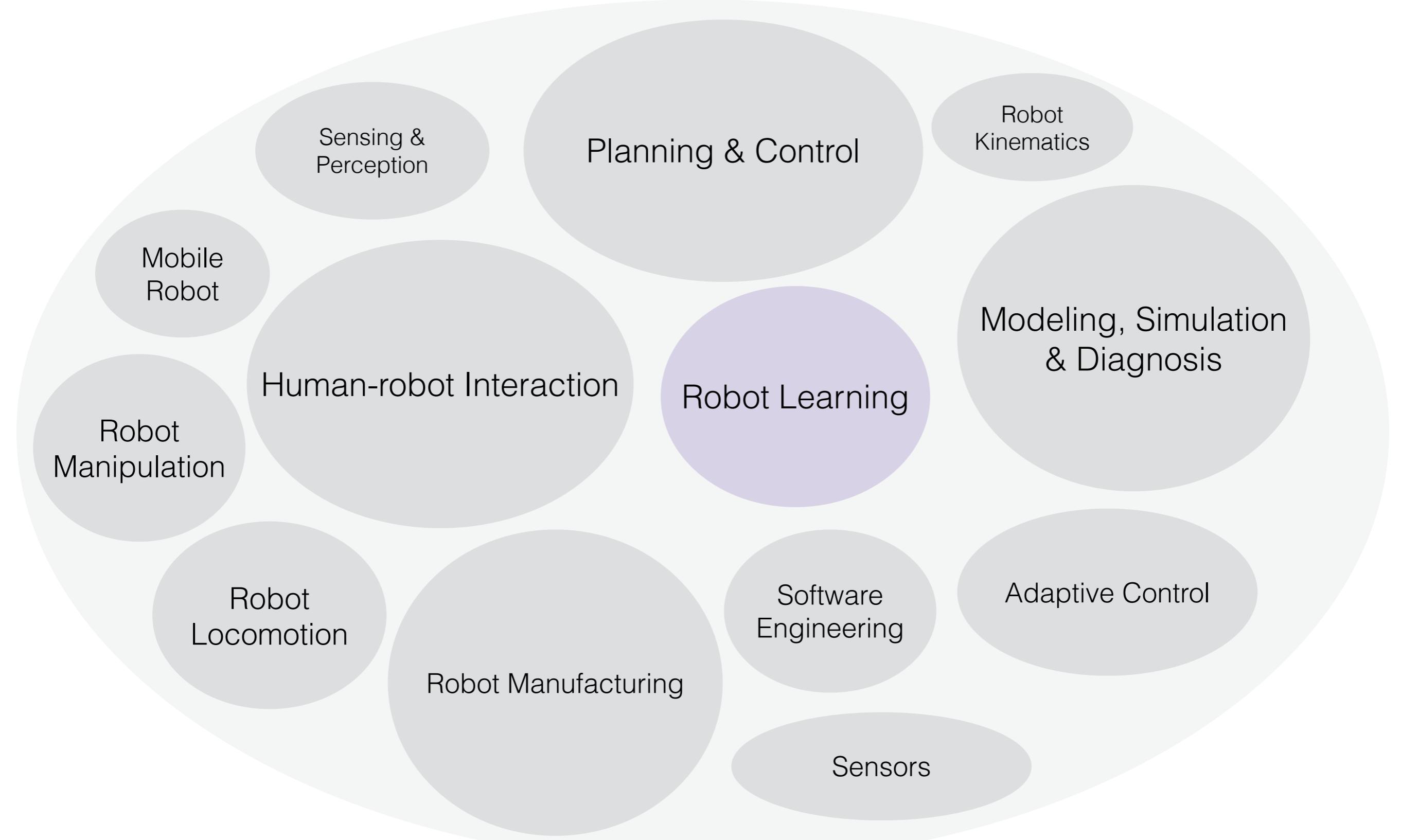
Semi-supervised Learning

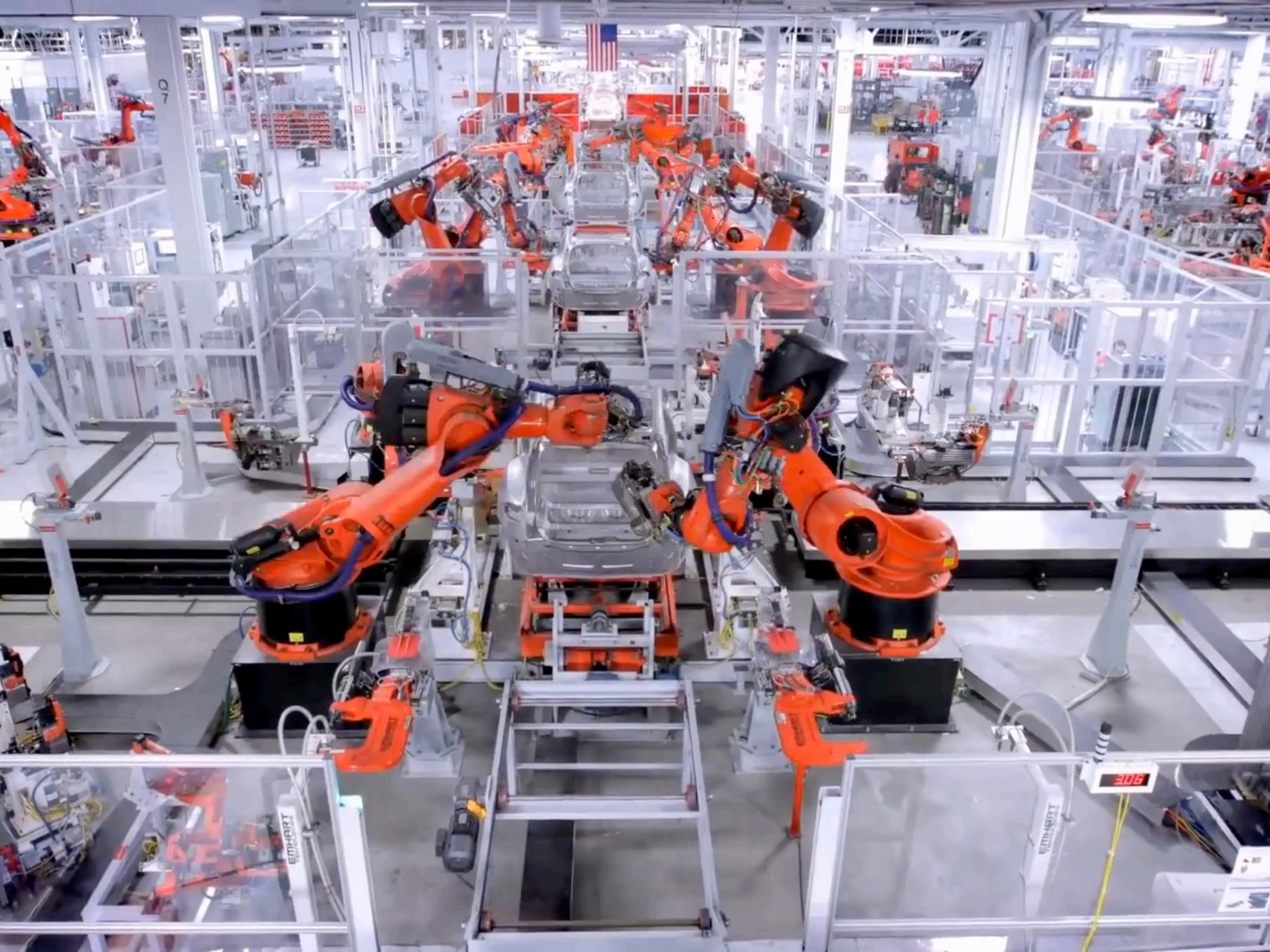
Self Learning

Feature Learning

Multi-Instance Learning

Robotics



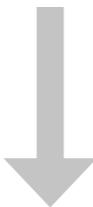


Robot Learning

Environment

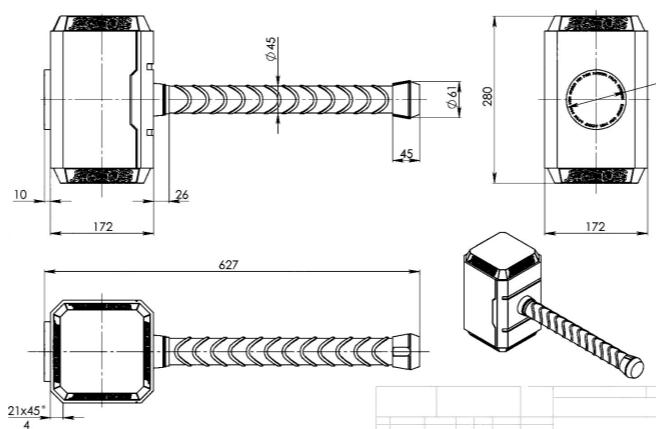


Structured

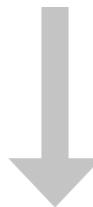


Unstructured

Object



Known



Unseen

Task



Pre-defined / Pre-programmed



Diverse and Novel

Supervised Learning

Image Classification

Model	image size	# parameters	Mult-Adds	Top 1 Acc. (%)	Top 5 Acc. (%)
Inception V2 [29]	224×224	11.2 M	1.94 B	74.8	92.2
NASNet-A (5 @ 1538)	299×299	10.9 M	2.35 B	78.6	94.2
Inception V3 [60]	299×299	23.8 M	5.72 B	78.8	94.4
Xception [9]	299×299	22.8 M	8.38 B	79.0	94.5
Inception ResNet V2 [58]	299×299	55.8 M	13.2 B	80.1	95.1
NASNet-A (7 @ 1920)	299×299	22.6 M	4.93 B	80.8	95.3
ResNeXt-101 (64 x 4d) [68]	320×320	83.6 M	31.5 B	80.9	95.6
PolyNet [69]	331×331	92 M	34.7 B	81.3	95.8
DPN-131 [8]	320×320	79.5 M	32.0 B	81.5	95.8
SENet [25]	320×320	145.8 M	42.3 B	82.7	96.2
NASNet-A (6 @ 4032)	331×331	88.9 M	23.8 B	82.7	96.2

Table 2. Performance of architecture search and other published state-of-the-art models ImageNet classification. Mult-Adds indicate the number of composite multiply-accumulate operations for a single image. Note that the composite multiple-accumulate operations are calculated for the image size reported in the table. Model size for [25] calculated from open-source implementation.

Zoph et al. Learning Transferable Architectures for Scalable Image Recognition

Instance Segmentation



Figure 5. More results of Mask R-CNN on COCO test images, using ResNet-101-FPN and running at 5 fps, with 35.7 mask AP (Table 1).

He et al. Mask R-CNN

Visual Question Answering

Method	VQA v2 test-dev			VQA v2 test-sd		
	All	Yes/no	Num	Other	All	Yes/no
Prior (most common answer in training set) [13]	-	-	-	-	25.08	41.20
LSTM Language only (blind model) [14]	-	-	-	-	44.08	67.01
Deeper LSTM Q score [17] as reported in [14]	-	-	-	-	54.26	73.46
MCB [13] as reported in [14]	-	-	-	-	62.27	78.82
UPMC-LPFG [1]	-	-	-	-	65.71	82.07
ESNUS	-	-	-	-	67.93	83.97
HDU-USyd-UNCC	-	-	-	-	66.77	81.89
Proposed model	-	-	-	-	68.09	84.50
ResNet features 7×7, single network	62.07	79.20	39.46	52.62	62.27	79.32
Image features from bottom-up attention, adaptive K, single network	65.32	81.82	44.21	56.05	65.67	82.20
ResNet features 7×7, ensemble	66.34	83.38	43.17	57.10	66.73	83.71
Image features from bottom-up attention, adaptive K, ensemble	69.87	86.08	45.99	68.60	70.34	86.64

Table 3. Comparison of our best model with competing methods. Excerpt from the official VQA v2 Leaderboard [1].

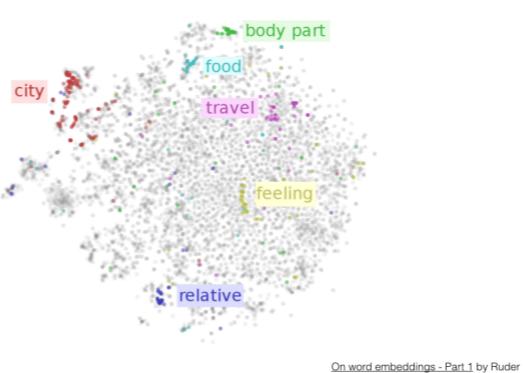
Teney et al. Tips and Tricks for Visual Question Answering: Learnings from the 2017 Challenge

Machine Translation

Source	"The reason Boeing are doing this is to cram more seats in to make their plane more competitive with our products," said Kevin Keniston, head of passenger comfort at Europe's Airbus.
PBMT	"La raison pour laquelle Boeing sont en train de faire, c'est de concentrer davantage de sièges pour prendre leur avion plus compétitif avec nos produits", a déclaré Kevin M. Keniston, chef du confort des passagers de l'Airbus de l'Europe.
GNMT	"La raison pour laquelle Boeing fait cela est de créer plus de sièges pour rendre son avion plus compétitif avec nos produits", a déclaré Kevin Keniston, chef du confort des passagers chez Airbus.
Human	"Boeing fait ça pour pouvoir caser plus de sièges et rendre ses avions plus compétitifs par rapport à nos produits", a déclaré Kevin Keniston, directeur de Confort Passager chez l'avionneur européen Airbus.
Source	When asked about this, an official of the American administration replied: "The United States is not conducting electronic surveillance aimed at offices of the World Bank and IMF in Washington."
PBMT	Interrogé à ce sujet, un responsable de l'administration américaine a répondu : "Les Etats-Unis n'est pas effectuer une surveillance électronique destiné aux bureaux de la Banque mondiale et du FMI à Washington".
GNMT	Interrogé à ce sujet, un fonctionnaire de l'administration américaine a répondu: "Les États-Unis n'effectuent pas de surveillance électronique à l'intention des bureaux de la Banque mondiale et du FMI à Washington".
Human	Interrogé sur le sujet, un responsable de l'administration américaine a répondu: "les Etats-Unis ne mènent pas de surveillance électronique visant les sièges de la Banque mondiale et du FMI à Washington".

Wu et al. Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation

Word Embeddings



On word embeddings - Part 1 by Ruder

Named Entity Recognition

contentSkip to site indexPoliticsSubscribeLogInSubscribeLogInToday's PaperAdvertisementSupported ORG byF.B.I. Agent Peter Strzok PERSON . Who Criticized Trump PERSON in Texts, Is FiredImagePeter Strzok, a top F.B.I. GPE counterintelligence agent who was taken off the special counsel investigation after his disparaging texts about President Trump PERSON were uncovered, was fired. CreditT.J. Kirkpatrick PERSON for The New York TimesBy Adam Goldman ORG and Michael S. SchmidtAug PERSON 13 CARDINAL , 2016WASHINGTON CARDINAL — Peter Strzok PERSON , the F.B.I. GPE senior counterintelligence agent who disparaged President Trump PERSON in inflammatory text messages and helped oversee the Hillary Clinton PERSON email and Russian GPE investigations, has been fired for violating bureau policies, Mr. Strzok PERSON , a lawyer said Monday DATE . Mr. Trump and his allies seized on the texts — exchanged during the 2016 DATE campaign with a former F.B.I. GPE lawyer, Lisa Page — in PERSON assailing the Russia GPE investigation as an illegitimate "witch hunt." Mr. Strzok PERSON , who rose over 20 years DATE at the F.B.I. GPE to become one of its most experienced counterintelligence agents, was a key figure in the early months DATE of the inquiry Along with writing the texts, Mr. Strzok PERSON was accused of sending a highly sensitive search warrant to his personal email account The F.B.I. GPE had been under immense political pressure by Mr. Trump PERSON to dismiss Mr. Strzok PERSON , who was removed last summer DATE from the staff of the special counsel, Robert S. Mueller III PERSON . The president has repeatedly denounced Mr. Strzok PERSON in posts on

Named Entity Recognition and Classification with Scikit-Learn by Susan Li Esteves et al. Named Entity Recognition in Twitter using Images and Text

Question Answering

System	Dev EM	Dev F1	Test EM	Test F1
Leaderboard (Oct 8th, 2018)				
Human	-	-	82.3	91.2
#1 Ensemble - nlnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	90.5
#1 Single - nlnet	-	-	83.5	90.1
#2 Single - QANet	-	-	82.5	89.3
Published				
BIDAF+ELMo (Single)	-	85.8	-	-
R.M. Reader (Single)	78.9	86.3	79.5	86.6
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours				
BERT _{BASE} (Single)	80.8	88.5	-	-
BERT _{LARGE} (Single)	84.1	90.9	-	-
BERT _{LARGE} (Ensemble)	85.8	91.8	-	-
BERT _{LARGE} (Sgl.+TriviaQnA)	84.2	91.5	85.1	91.8
BERT _{LARGE} (Ens.+TriviaQnA)	86.2	92.2	87.4	93.2

Table 2: SQuAD results. The BERT ensemble is 7x systems which use different pre-training checkpoints and fine-tuning seeds.

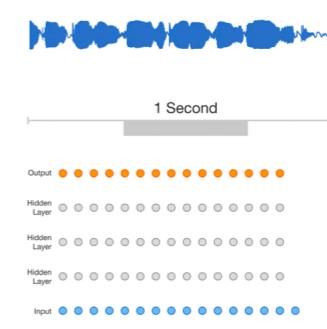
Devlin et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Speech Recognition

Senone set	Model/combination step	WER		WER	
		devset	test	ngram-LM	LSTM
9k	BLSTM	11.5	8.3	9.2	6.3
27k	BLSTM	11.4	8.0	9.2	6.3
9k-puhpum	BLSTM	11.3	8.0	9.2	5.4
9k-puhpum	BLSTM+ResNet+LACE+CNN-BLSTM	9.6	7.2	7.7	5.4
9k-puhpum	BLSTM+ResNet+LACE	9.7	7.4	7.8	5.4
27k	BLSTM+ResNet+LACE+CNN-BLSTM	9.7	7.3	7.8	5.5
-	BLSTM+ResNet+LACE	10.0	7.5	8.0	5.8
Confusion network combination		7.4		5.2	
+ LSTM rescoring		7.3		5.2	
+ ngram rescoring		7.2		5.2	
+ backchannel penalty		7.2		5.1	

Xiong et al. The Microsoft 2017 Conversational Speech Recognition System

Speech Synthesis (text-to-speech)



Van Den Oord et al. WaveNet: A Generative Model for Raw Audio

Supervised Learning

Image Classification

Model	image size	# parameters	Mult-Adds	Top 1 Acc. (%)	Top 5 Acc. (%)
Inception V2 [29]	224×224	11.2 M	1.94 B	74.8	92.2
NASNet-A (5 @ 1538)	299×299	10.9 M	2.35 B	78.6	94.2
Inception V3 [60]	299×299	23.8 M	5.72 B	78.8	94.4
Xception [9]	299×299	22.8 M	8.38 B	79.0	94.5
Inception ResNet V2 [58]	299×299	55.8 M	13.2 B	80.1	95.1
NASNet-A (7 @ 1920)	299×299	22.6 M	4.93 B	80.8	95.3
ResNeXt-101 (64 × 4d) [68]	320×320	83.6 M	31.5 B	80.9	95.6
PolyNet [69]	331×331	92 M	34.7 B	81.3	95.5
DPN-131 [8]	320×320	79.5 M	32.0 B	81.5	95.5
SENet [25]	320×320	145.8 M	42.3 B	82.7	95.6
NASNet-A (6 @ 4032)	331×331	88.9 M	23.8 B	82.7	95.2

Table 2. Performance of architecture search and other published state-of-the-art models on ImageNet classification. Mult-Adds indicate the number of composite multiply-accumulate operations for a single image. Note that the composite multiple-accumulate operations are calculated for the image size reported in the table. Model size for [25] calculated from open-source implementation.

Instance Segmentation



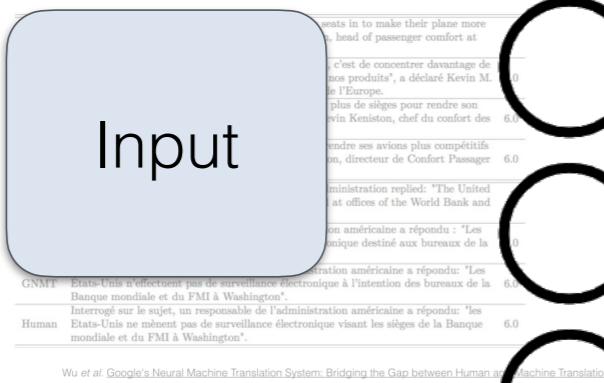
Zoph et al. Learning Transferable Architectures for Scalable Image Recognition

Visual Question Answering

Method	VQA v2 test-dev			VQA v2 test-old				
	All	Yes/no	Numth.	Other	All	Yes/no	Numth.	Other
Prior (most common answer in training set) [1]	-	-	-	-	25.08	18.08	3.17	3.17
LSTM Language only (blind model) [1]	-	-	-	-	44.26	47.40	31.55	27.37
Deeper LSTM Q & vnn. [17] as reported in [1]	-	-	-	-	54.22	73.40	35.18	41.83
MCB [1] as reported in [1]	-	-	-	-	62.27	78.42	38.28	53.36
UPMC-LCPI [1]	-	-	-	-	65.71	82.07	41.06	57.12
RCNN [1]	-	-	-	-	67.29	82.07	41.06	57.17
UNICORN	-	-	-	-	68.77	81.80	46.29	58.30
HDI-USyd-UNCC	-	-	-	-	68.09	84.50	45.39	59.01
Proposed model	-	-	-	-	62.07	79.20	39.46	52.62
ResNet features T×7, single network	62.07	79.20	39.46	52.62	62.27	79.32	39.77	52.59
Image features from bottom-up attention, adaptive K , single network	65.32	81.82	44.21	56.05	65.67	82.20	43.90	56.26
ResNet features T×7, ensemble	66.34	83.38	43.19	57.40	66.73	83.71	43.77	57.20
Image features from bottom-up attention, adaptive K , ensemble	69.87	86.08	48.99	68.60	70.34	86.64	48.64	61.15

Table 3. Comparison of our best model with competing methods. Excerpt from the official VQA v2 Leaderboard [1].

Machine Translation



Question Answering

System	Dev EM	Dev F1	Test EM	Test F1
Leaderboard (Oct 8th, 2018)				
Human	-	-	82.4	91.2
#1 Ensemble - nfinet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	91.5
#1 Single - nfinet	-	-	83.5	90.3
#2 Single - QANet	-	-	82.5	89.3
Published				
BIDAF+ELMo (Single)	-	85.8	-	-
R.M. Reader (Single)	78.9	86.3	79.5	86.6
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours	-	-	-	-
BERT _{BASE} (Single)	80.8	88.5	-	-
BERT _{LARGE} (Single)	84.1	90.9	-	-
BERT _{LARGE} (Ensemble)	85.8	91.8	-	-
BERT _{LARGE} (Sgl.+TriviaQnA)	84.2	91.5	85.1	91.8
BERT _{LARGE} (Ens.+TriviaQnA)	86.2	92.2	87.4	93.2

Table 2: SQuAD results. The BERT ensemble is 7x systems which use different pre-training checkpoints and fine-tuning seeds.

Devlin et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

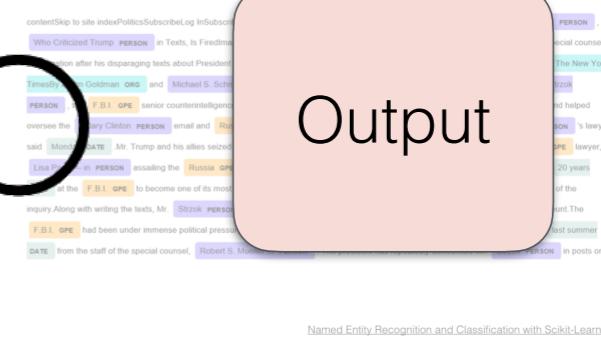
Word Embeddings

Speech Recognition

Senone set	Model/combination step	WER		WER	
		devset	test	devset	test
9k	BLSTM	11.5	8.3	9.2	6.3
27k	BLSTM	11.4	8.0	9.2	6.3
9k	BLSTM	11.3	8.0	9.2	6.3
9k-puhpum	BLSTM+ResNet+LACE+CNN-BLSTM	9.6	7.2	7.7	5.4
9k-puhpum	BLSTM+ResNet+LACE	9.7	7.4	7.8	5.4
27k	BLSTM+ResNet+LACE	9.7	7.3	7.8	5.5
9k	Confusion network combination	-	-	7.4	5.2
-	+ LSTM rescoring	-	-	7.3	5.2
-	+ ngram rescoring	-	-	7.2	5.2
-	+ backchannel penalty	-	-	7.2	5.1

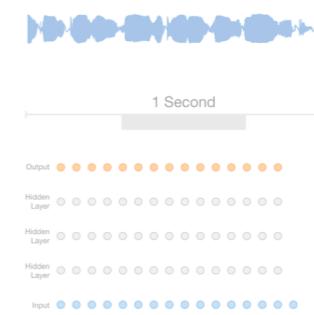
Xiong et al. The Microsoft 2017 Conversational Speech Recognition System

Named Entity Recognition



Named Entity Recognition and Classification with Scikit-Learn by Susan Li
Esteves et al. Named Entity Recognition in Twitter using Images and Text

Speech Synthesis (text-to-speech)



Van Den Oord et al. WaveNet: A Generative Model for Raw Audio

Image Classification

Image Classification

Model	image size	# parameters	Mult-Adds	Top 1 Acc. (%)	Top 5 Acc. (%)
Inception V2 [29]	224×224	11.2 M	1.94 B	74.8	92.2
NASNet-A (5 @ 1538)	299×299	10.9 M	2.35 B	78.6	94.2
Inception V3 [60]	299×299	23.8 M	5.72 B	78.8	94.4
Xception [9]	299×299	22.8 M	8.38 B	79.0	94.5
Inception ResNet V2 [58]	299×299	55.8 M	13.2 B	80.1	95.1
NASNet-A (7 @ 1920)	299×299	22.6 M	4.93 B	80.8	95.3
ResNeXt-101 (64 × 4d) [68]	320×320	83.6 M	31.5 B	80.9	95.6
PolyNet [69]	331×331	92 M	34.7 B	81.3	95.5
DPN-131 [8]	320×320	79.5 M	32.0 B	81.5	95.5
SENet [25]	320×320	145.8 M	42.3 B	82.7	95.6
NASNet-A (6 @ 4032)	331×331	88.9 M	23.8 B	82.7	95.2

Table 2. Performance of architecture search and other published state-of-the-art models. ImageNet classification. Mult-Adds indicate the number of composite multiply-accumulate operations for a single image. Note that the composite multiple-accumulate operations are calculated for the image size reported in the [1]. Model size from [25] calculated from open-source implementation.

Instance Segmentation



Zoph et al. Learning Transferable Architectures for Scalable Image Recognition

Visual Question Answering

Method	VQA v2 test-dev			VQA v2 test-old				
	All	Yes/no	Numth.	Other	All	Yes/no	Numth.	Other
Prior most common answer in training set [1]	-	-	-	-	25.08	47.01	31.55	3.17
LSTM Language only (blind model) [1]	-	-	-	-	44.26	67.40	31.55	27.37
Deeper LSTM Q & vnn. [17] as reported in [1]	-	-	-	-	54.22	73.40	35.18	41.83
MCB [1] as reported in [1]	-	-	-	-	62.27	78.42	38.28	53.36
UPMC-LAPB [1]	-	-	-	-	65.71	82.07	41.06	57.12
Proposed model	-	-	-	-	67.29	83.44	41.77	57.97
HDI-USyd-UNCC	-	-	-	-	68.77	81.80	46.29	58.30
Proposed model	-	-	-	-	68.09	84.50	45.39	59.01
ResNet features T×7, single network	62.07	79.20	39.46	52.62	62.27	79.32	39.77	52.59
Image features from bottom-up attention, adaptive K , single network	65.32	81.82	44.21	56.05	65.67	82.20	43.90	56.26
ResNet features T×7, ensemble	66.34	83.58	43.19	57.40	66.73	83.71	43.77	57.20
Image features from bottom-up attention, adaptive K , ensemble	69.87	86.08	48.99	68.03	70.34	86.64	48.64	61.15

Table 3. Comparison of our best model with competing methods. Excerpt from the official VQA v2 Leaderboard [1].

Image



Question Answering

- Input Question:
Where do water droplets collide with ice crystals to form precipitation?
- Input Paragraph:
... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. ...
- Output Answer:
within a cloud

System	Dev EM	Dev F1	Test EM	Test F1
Leaderboard (Oct 8th, 2018)				
Human	-	-	82.1	91.2
#1 Ensemble - nlnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	91.5
#1 Single - nlnet	-	-	83.5	90.3
#2 Single - QANet	-	-	82.5	89.3
Published				
BIDAF+ELMo (Single)	-	85.8	-	-
R.M. Reader (Single)	78.9	86.3	79.5	86.6
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours	-	-	-	-
BERT _{BASE} (Single)	80.8	88.5	-	-
BERT _{LARGE} (Single)	84.1	90.9	-	-
BERT _{LARGE} (Ensemble)	85.8	91.8	-	-
BERT _{LARGE} (Sgl.+TriviaQA)	84.2	91.1	85.1	91.8
BERT _{LARGE} (Ens.+TriviaQA)	86.2	92.2	87.4	93.2

Table 2: SQuAD results. The BERT ensemble is 7x systems which use different pre-training checkpoints and fine-tuning seeds.

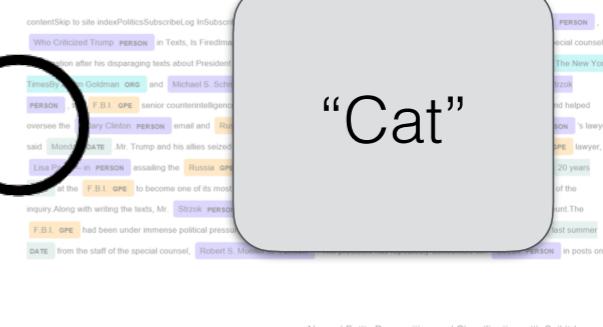
Devlin et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Speech Recognition

Senone set	Model/combination step	WER		WER	
		devset	test	ngram-LM	LSTM-LMs
9k	BLSTM	11.5	8.3	9.2	6.3
27k	BLSTM	11.4	8.0	9.2	6.3
9k-puhpum	BLSTM	11.3	8.0	9.2	6.3
9k-puhpum	BLSTM+ResNet+LACE+CNN-BLSTM	9.6	7.2	7.7	5.4
9k-puhpum	BLSTM+ResNet+LACE	9.7	7.4	7.8	5.4
27k	BLSTM+ResNet+LACE	9.7	7.3	7.8	5.5
	BLSTM+ResNet+LACE	10.0	7.5	8.0	5.8
	Confusion network combination			7.4	5.2
	+ LSTM rescoring			7.3	5.2
	+ ngram rescoring			7.2	5.2
	+ backchannel penalty			7.2	5.1

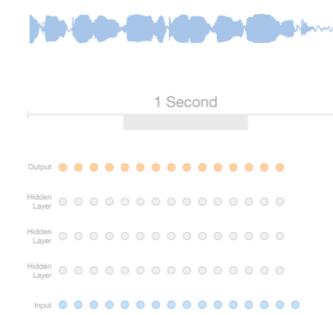
Xiong et al. The Microsoft 2017 Conversational Speech Recognition System

Named Entity Recognition



Named Entity Recognition and Classification with Scikit-Learn by Susan Li Esteves et al. Named Entity Recognition in Twitter using Images and Text

Speech Synthesis (text-to-speech)



Van Den Oord et al. WaveNet: A Generative Model for Raw Audio

Machine Translation

Image Classification

Model	image size	# parameters	Mult-Adds	Top 1 Acc. (%)	Top 5 Acc. (%)
Inception V2 [29]	224×224	11.2 M	1.94 B	74.8	92.2
NASNet-A (5 @ 1538)	299×299	10.9 M	2.35 B	78.6	94.2
Inception V3 [60]	299×299	23.8 M	5.72 B	78.8	94.4
Xception [9]	299×299	22.8 M	8.38 B	79.0	94.5
Inception ResNet V2 [58]	299×299	55.8 M	13.2 B	80.1	95.1
NASNet-A (7 @ 1920)	299×299	22.6 M	4.93 B	80.8	95.3
ResNeXt-101 (64 × 4d) [68]	320×320	83.6 M	31.5 B	80.9	95.6
PolyNet [69]	331×331	92 M	34.7 B	81.3	95.5
DPN-131 [8]	320×320	79.5 M	32.0 B	81.5	95.5
SENet [25]	320×320	145.8 M	42.3 B	82.7	95.6
NASNet-A (6 @ 4032)	331×331	88.9 M	23.8 B	82.7	95.2

Table 2. Performance of architecture search and other published state-of-the-art models. ImageNet classification. Mult-Adds indicate the number of composite multiply-accumulate operations for a single image. Note that the composite multiple-accumulate operations are calculated for the image size reported in the table. Model size for [25] calculated from open-source implementation.

Instance Segmentation



Zoph et al. Learning Transferable Architectures for Scalable Image Recognition

Visual Question Answering

Method	VQA v2 test-dev			VQA v2 test-old				
	All	Yes/no	Numth.	Other	All	Yes/no	Numth.	Other
Prior (most common answer in training set) [1]	-	-	-	-	25.08	16.70	3.17	3.17
LSTM Language only (blind model) [1]	-	-	-	-	44.26	47.40	31.55	27.37
Deeper LSTM Q & vnn. [17] as reported in [1]	-	-	-	-	54.22	73.40	35.18	41.83
MCB [1] as reported in [1]	-	-	-	-	62.27	78.42	38.28	53.36
UPMC-LPBP [1]	-	-	-	-	65.71	82.07	41.06	57.12
“+” [1]	-	-	-	-	67.29	84.20	43.77	59.97
LiONUS	-	-	-	-	68.77	81.80	46.29	58.30
HDI-USyd-UNCC	-	-	-	-	68.09	84.50	45.39	59.01
Proposed model	-	-	-	-	62.07	79.20	39.46	52.62
ResNet features T×7, single network	62.07	79.20	39.46	52.62	62.27	79.32	39.77	52.59
Image features from bottom-up attention, adaptive K , single network	65.32	81.82	44.21	56.05	65.67	82.20	43.90	56.26
ResNet features T×7, ensemble	66.34	83.38	43.10	57.40	66.73	83.71	43.77	57.20
Image features from bottom-up attention, adaptive K , ensemble	69.87	86.08	48.99	68.60	70.34	86.64	51.15	

Table 3. Comparison of our best model with competing methods. Excerpt from the official VQA v2 Leaderboard [1].

English sentence

France is never cold
in September

GNMT : Etats-Unis n'effectuent pas de surveillance électronique à l'intention des bureaux de la Banque mondiale et du FMI à Washington.
Interrogé sur le sujet, un responsable de l'administration américaine a répondu: "les Etats-Unis ne mènent pas de surveillance électronique visant les sièges de la Banque mondiale et du FMI à Washington".

Wu et al. Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation



French sentence

la france est jamais
froid en septembre

contentSkip to site indexPoliticsSub
Who Criticized Trump PERSON
Teleology in Goldblum GPE
oversee the
said Movistar DATE Mr. Trump
Lisa DiCarlo PERSON
assault
inquiry Along with writing the texts
FBI GPE had been under investigation
DATE from the staff of the special counsel

Named Entity Recognition and Classification with Scikit-Learn by Susan Li
Esteves et al. Named Entity Recognition in Twitter using Images and Text

Question Answering

System	Dev EM	Dev F1	Test EM	Test F1
Leaderboard (Oct 8th, 2018)				
Human	-	-	82.1	91.2
#1 Ensemble - nlnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	91.5
#1 Single - nlnet	-	-	83.5	90.8
#2 Single - QANet	-	-	82.5	89.3
Published				
BIDAF+ELMo (Single)	-	85.8	-	-
R.M. Reader (Single)	78.9	86.3	79.5	86.6
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours				
BERT _{BASE} (Single)	80.8	88.5	-	-
BERT _{LARGE} (Single)	84.1	90.9	-	-
BERT _{LARGE} (Ensemble)	85.8	91.8	-	-
BERT _{LARGE} (Sgl.+TriviaQnA)	84.2	91.5	85.1	91.8
BERT _{LARGE} (Ens.+TriviaQnA)	86.2	92.2	87.4	93.2

Table 2: SQuAD results. The BERT ensemble is 7x systems which use different pre-training checkpoints and fine-tuning seeds.

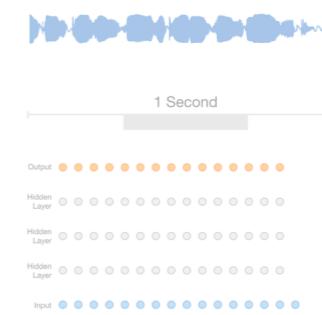
Devlin et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Speech Recognition

Senone set	Model/combination step	WER		WER	
		devset	test	ngram-LM	devset
9k	BLSTM	11.5	8.3	9.2	6.3
27k	BLSTM	11.4	8.0	9.2	6.3
9k	BLSTM	11.3	8.0	9.2	6.3
9k-puhpum	BLSTM+ResNet+LACE+CNN-BLSTM	9.6	7.2	7.7	5.4
9k-puhpum	BLSTM+ResNet+LACE	9.7	7.4	7.8	5.4
27k	BLSTM+ResNet+LACE	9.7	7.3	7.8	5.5
9k	Confusion network combination	-	-	7.4	5.2
-	+ LSTM rescoring	-	-	7.3	5.2
-	+ ngram rescoring	-	-	7.2	5.2
-	+ backchannel penalty	-	-	7.2	5.1

Xiong et al. The Microsoft 2017 Conversational Speech Recognition System

Speech Synthesis (text-to-speech)



Van Den Oord et al. WaveNet: A Generative Model for Raw Audio

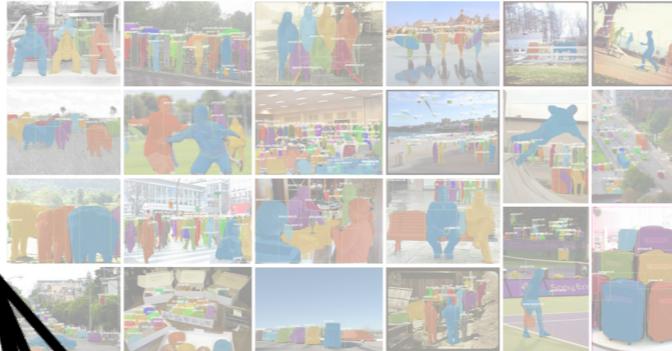
Automatic Speech Recognition

Image Classification

Model	image size	# parameters	Mult-Adds	Top 1 Acc. (%)	Top 5 Acc. (%)
Inception V2 [29]	224×224	11.2 M	1.94 B	74.8	92.2
NASNet-A (5 @ 1538)	299×299	10.9 M	2.35 B	78.6	94.2
Inception V3 [60]	299×299	23.8 M	5.72 B	78.8	94.4
Xception [9]	299×299	22.8 M	8.38 B	79.0	94.5
Inception ResNet V2 [58]	299×299	55.8 M	13.2 B	80.1	95.1
NASNet-A (7 @ 1920)	299×299	22.6 M	4.93 B	80.8	95.3
ResNeXt-101 (64 × 4d) [68]	320×320	83.6 M	31.5 B	80.9	95.6
PolyNet [69]	331×331	92 M	34.7 B	81.3	95.5
DPN-131 [8]	320×320	79.5 M	32.0 B	81.5	95.5
SENet [25]	320×320	145.8 M	42.3 B	82.7	95.6
NASNet-A (6 @ 4032)	331×331	88.9 M	23.8 B	82.7	95.2

Table 2. Performance of architecture search and other published state-of-the-art models on ImageNet classification. Mult-Adds indicate the number of composite multiply-accumulate operations for a single image. Note that the composite multiple-accumulate operations are calculated for the image size reported in the table. Model size for [25] calculated from open-source implementation.

Instance Segmentation



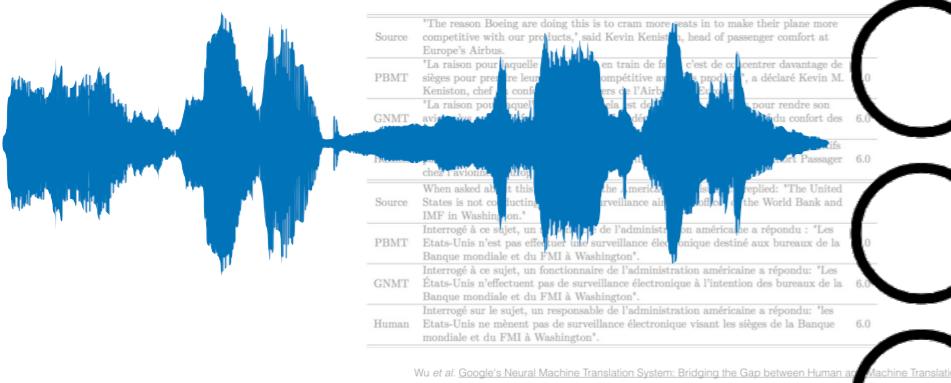
Zoph et al. Learning Transferable Architectures for Scalable Image Recognition

Visual Question Answering

Method	VQA v2 test-dev			VQA v2 test-old				
	All	Yes/no	Numth.	Other	All	Yes/no	Numth.	Other
Prior most common answer in training set [1]	-	-	-	-	25.08	19.08	3.17	3.17
LSTM Language only (blind model) [1]	-	-	-	-	44.26	47.40	31.55	27.37
Deeper LSTM Q score [17] as reported in [1]	-	-	-	-	54.22	73.40	35.18	41.83
MCB [10] as reported in [1]	-	-	-	-	62.27	78.42	38.28	53.36
UPMC-LCPI [1]	-	-	-	-	65.71	82.07	41.06	57.12
UCON [18]	-	-	-	-	67.27	84.00	43.77	59.97
HCU-USyd-UNCC	-	-	-	-	68.77	81.80	46.29	58.30
Proposed model	-	-	-	-	68.09	84.50	45.39	59.01
ResNet features T×7, single network	62.07	79.20	39.46	52.62	62.27	79.32	39.77	52.59
Image features from bottom-up attention, adaptive K , single network	65.32	81.82	44.21	56.05	65.67	82.20	43.90	56.26
ResNet features T×7, ensemble	66.34	83.38	43.19	57.40	66.73	83.71	43.77	57.20
Image features from bottom-up attention, adaptive K , ensemble	69.87	86.08	48.99	68.64	70.34	86.60	48.64	61.15

Table 3. Comparison of our best model with competing methods. Excerpt from the official VQA v2 Leaderboard [1].

Waveform Machine Translation



Wu et al. Google’s Neural Machine Translation System: Bridging the Gap between Human and Machine Translation

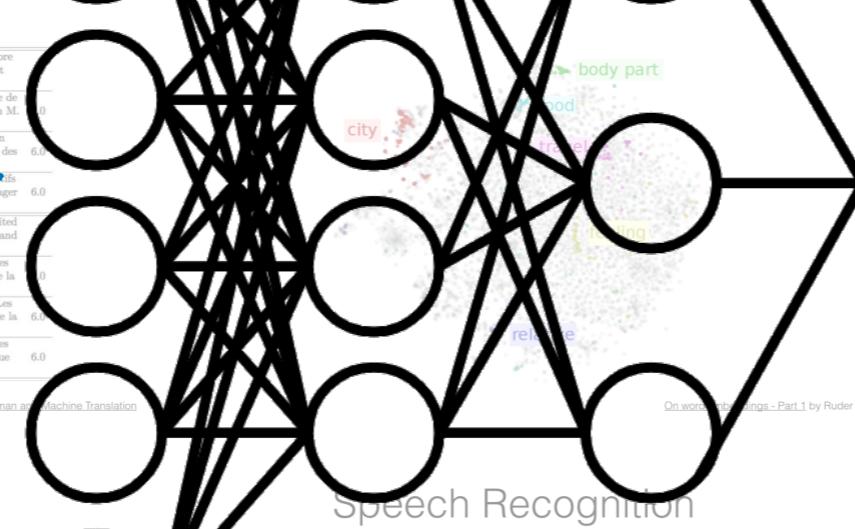
Question Answering

System	Dev EM	Dev F1	Test EM	Test F1
Leaderboard (Oct 8th, 2018)				
Human	-	-	82.9	91.2
#1 Ensemble - nifnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	91.5
#1 Single - nifnet	-	-	83.5	90.3
#2 Single - QANet	-	-	82.5	89.3
Published				
BIDAF+ELMo (Single)	-	85.8	-	-
R.M. Reader (Single)	78.9	86.3	79.5	86.6
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours				
BERT _{BASE} (Single)	80.8	88.5	-	-
BERT _{LARGE} (Single)	84.1	90.9	-	-
BERT _{LARGE} (Ensemble)	85.8	91.8	-	-
BERT _{LARGE} (Sgl.+TriviaQnA)	84.2	91.5	85.1	91.8
BERT _{LARGE} (Ens.+TriviaQnA)	86.2	92.2	87.4	93.2

Table 2: SQuAD results. The BERT ensemble is 7x systems which use different pre-training checkpoints and fine-tuning seeds.

Devlin et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Word Embeddings



On word embeddings - Part 1 by Ruder

Text Named Entity Recognition

This is a supervised learning method

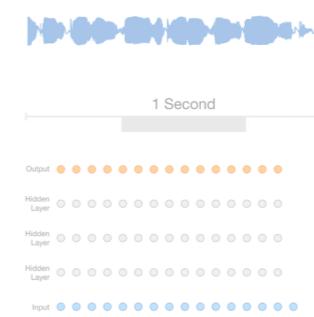
Named Entity Recognition and Classification with Scikit-Learn by Susan Li
Esteves et al. Named Entity Recognition in Twitter using Images and Text

Speech Recognition

Senone set	Model/combination step	WER		WER	
		devset	test	ngram-LM	LSTM-LMs
9k	BLSTM	11.5	8.3	9.2	6.3
27k	BLSTM	11.4	8.0	9.2	6.3
9k-puhpum	BLSTM	11.3	8.0	9.2	5.4
9k-puhpum	BLSTM+ResNet+LACE+CNN-BLSTM	9.6	7.2	7.7	5.4
9k-puhpum	BLSTM+ResNet+LACE	9.7	7.4	7.8	5.4
27k	BLSTM+ResNet+LACE+CNN-BLSTM	9.7	7.3	7.8	5.5
27k	BLSTM+ResNet+LACE	10.0	7.5	8.0	5.8
Confusion network combination		7.4		5.2	
+ LSTM rescoring		7.3		5.2	
+ ngram rescoring		7.2		5.2	
+ backchannel penalty		7.2		5.1	

Xiong et al. The Microsoft 2017 Conversational Speech Recognition System

Speech Synthesis (text-to-speech)



Van Den Oord et al. WaveNet: A Generative Model for Raw Audio

Supervised Learning

Image Classification

Model	image size	# parameters	Mult-Adds	Top 1 Acc. (%)	Top 5 Acc. (%)
Inception V2 [29]	224×224	11.2 M	1.94 B	74.8	92.2
NASNet-A (5 @ 1538)	299×299	10.9 M	2.35 B	78.6	94.2
Inception V3 [60]	299×299	23.8 M	5.72 B	78.8	94.4
Xception [9]	299×299	22.8 M	8.38 B	79.0	94.5
Inception ResNet V2 [58]	299×299	55.8 M	13.2 B	80.1	95.1
NASNet-A (7 @ 1920)	299×299	22.6 M	4.93 B	80.8	95.3
ResNeXt-101 (64 × 4d) [68]	320×320	83.6 M	31.5 B	80.9	95.6
PolyNet [69]	331×331	92 M	34.7 B	81.3	95.5
DPN-131 [8]	320×320	79.5 M	32.0 B	81.5	95.5
SENet [25]	320×320	145.8 M	42.3 B	82.7	95.6
NASNet-A (6 @ 4032)	331×331	88.9 M	23.8 B	82.7	95.2

Table 2. Performance of architecture search and other published state-of-the-art models on ImageNet classification. Mult-Adds indicate the number of composite multiply-accumulate operations for a single image. Note that the composite multiple-accumulate operations are calculated for the image size reported in the table. Model size for [25] calculated from open-source implementation.

Instance Segmentation



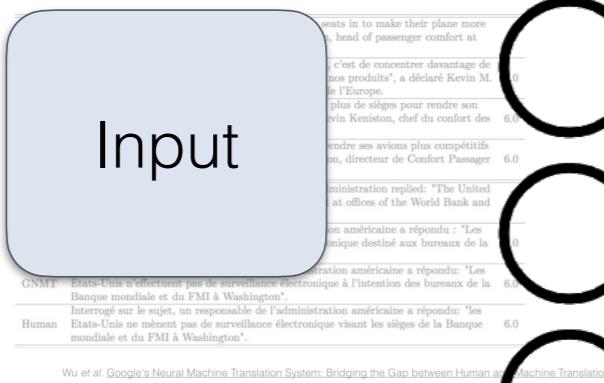
Zoph et al. Learning Transferable Architectures for Scalable Image Recognition

Visual Question Answering

Method	VQA v2 test-dev			VQA v2 test-old				
	All	Yes/no	Numth.	Other	All	Yes/no	Numth.	Other
Prior (most common answer in training set) [1]	-	-	-	-	25.08	18.08	3.17	3.17
LSTM Language only (blind model) [1]	-	-	-	-	44.26	47.40	31.55	27.37
Deeper LSTM Q & vnn. [17] as reported in [1]	-	-	-	-	54.22	73.40	35.18	41.83
MCB [1] as reported in [1]	-	-	-	-	62.27	78.42	38.28	53.36
UPMC-LCPI [1]	-	-	-	-	65.71	82.07	41.06	57.12
RCNN [1]	-	-	-	-	67.29	83.40	43.77	59.97
UNICOS	-	-	-	-	68.77	81.80	46.29	58.30
HDI-USyd-UNCC	-	-	-	-	68.09	84.50	45.39	59.01
Proposed model	-	-	-	-	62.07	79.20	39.46	52.62
ResNet features T×7, single network	62.07	79.20	39.46	52.62	62.27	79.32	39.77	52.59
Image features from bottom-up attention, adaptive K , single network	65.32	81.82	44.21	56.05	65.67	82.20	43.90	56.26
ResNet features T×7, ensemble	66.34	83.38	43.19	57.40	66.73	83.71	43.77	57.20
Image features from bottom-up attention, adaptive K , ensemble	69.87	86.08	48.99	68.60	70.34	86.64	48.64	61.15

Table 3. Comparison of our best model with competing methods. Excerpt from the official VQA v2 Leaderboard [1].

Machine Translation



Wu et al. Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation

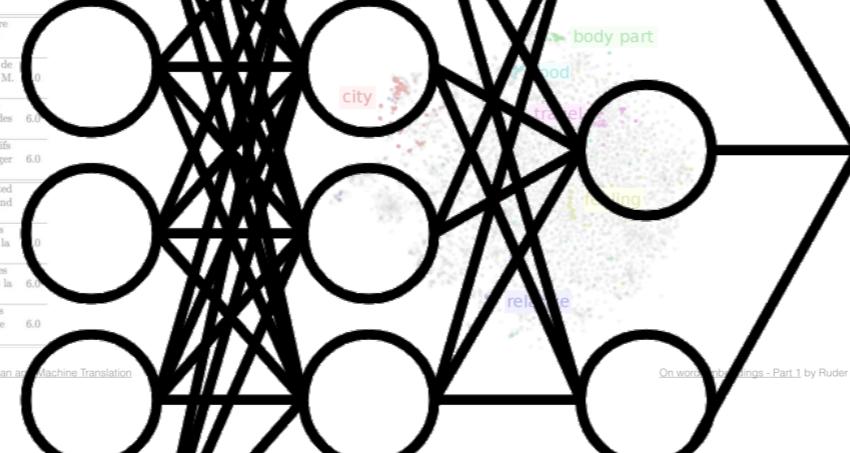
Question Answering

System	Dev EM	Dev F1	Test EM	Test F1
Leaderboard (Oct 8th, 2018)				
Human	-	-	82.4	91.2
#1 Ensemble - nInet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	91.5
#1 Single - nInet	-	-	83.5	90.3
#2 Single - QANet	-	-	82.5	89.3
Published				
BIDAF+ELMo (Single)	-	85.8	-	-
R.M. Reader (Single)	78.9	86.3	79.5	86.6
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours	-	-	-	-
BERT _{BASE} (Single)	80.8	88.5	-	-
BERT _{LARGE} (Single)	84.1	90.9	-	-
BERT _{LARGE} (Ensemble)	85.8	91.8	-	-
BERT _{LARGE} (Sgl.+TriviaQnA)	84.2	91.5	85.1	91.8
BERT _{LARGE} (Ens.+TriviaQnA)	86.2	92.2	87.4	93.2

Table 2: SQuAD results. The BERT ensemble is 7x systems which use different pre-training checkpoints and fine-tuning seeds.

Devlin et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Word Embeddings

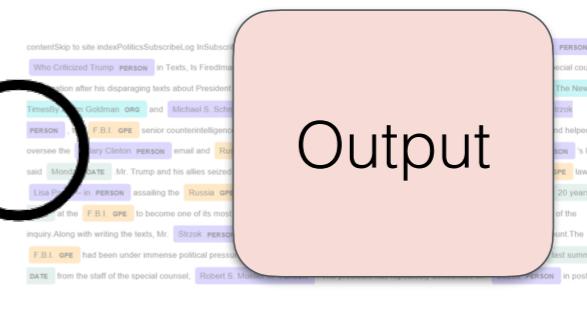


On word embeddings - Part 1 by Ruder

Senone set	Model/combination step	WER		WER	
		devset	test	ngram-LM	LSTM-LMs
9k	BLSTM	11.5	8.3	9.2	6.3
27k	BLSTM	11.4	8.0	9.2	6.3
9k	BLSTM	11.3	8.0	9.2	6.3
9k-puhpum	BLSTM+ResNet+LACE+CNN-BLSTM	9.6	7.2	7.7	5.4
9k-puhpum	BLSTM+ResNet+LACE	9.7	7.4	7.8	5.4
27k	BLSTM+ResNet+LACE	9.7	7.3	7.8	5.5
9k	Confusion network combination	-	-	7.4	5.2
-	+ LSTM rescoring	-	-	7.3	5.2
-	+ ngram rescoring	-	-	7.2	5.2
-	+ backchannel penalty	-	-	7.2	5.1

Xiong et al. The Microsoft 2017 Conversational Speech Recognition System

Named Entity Recognition



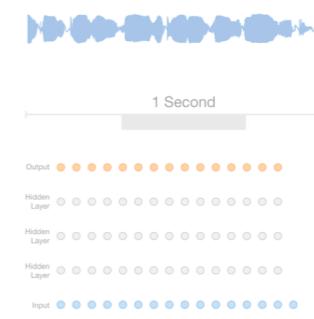
Named Entity Recognition and Classification with Scikit-Learn by Susan Li
Esteves et al. Named Entity Recognition in Twitter using Images and Text

Speech Recognition

Word Error Rate	Model/combination step	WER		WER	
		devset	test	devset	test
9k	BLSTM	11.5	8.3	9.2	6.3
27k	BLSTM	11.4	8.0	9.2	6.3
9k	BLSTM	11.3	8.0	9.2	6.3
9k-puhpum	BLSTM+ResNet+LACE+CNN-BLSTM	9.6	7.2	7.7	5.4
9k-puhpum	BLSTM+ResNet+LACE	9.7	7.4	7.8	5.4
27k	BLSTM+ResNet+LACE	9.7	7.3	7.8	5.5
9k	Confusion network combination	-	-	7.4	5.2
-	+ LSTM rescoring	-	-	7.3	5.2
-	+ ngram rescoring	-	-	7.2	5.2
-	+ backchannel penalty	-	-	7.2	5.1

Xiong et al. The Microsoft 2017 Conversational Speech Recognition System

Speech Synthesis (text-to-speech)



Van Den Oord et al. WaveNet: A Generative Model for Raw Audio

Supervised Learning

Image Classification

Model	image size	# parameters	Mult-Adds	Top 1 Acc. (%)	Top 5 Acc. (%)
Inception V2 [70]	299×299	11.5M	1.64B	74.8	92.2
NASNet-A (5@ 1538)	299×299	6.9M	2.35 B	78.6	94.2
Inception V3 [60]	299×299	11.5M	5.72 B	78.8	94.4
Xception [9]	299×299	11.5M	8.38 B	79.0	94.5
Inception ResNet V2 [58]	299×299	11.5M	13.2 B	80.1	95.1
NASNet-A (7@ 1920)	299×299	11.5M	4.93 B	80.8	95.3
MobileNet v2 [61]	320×320	83.6M	2.02 B	80.9	95.6
PolyNet [69]	333×333	92 M	34.7 B	81.3	95.5
DPN-131 [8]	320×320	79.5M	32.0 B	81.5	95.5
SENet [25]	320×320	145.8M	42.3 B	82.7	95.6
NASNet-A (6@ 4032)	333×333	88.9 M	23.8 B	82.7	95.2

Table 1. Performance of architecture search and other published state-of-the-art models on ImageNet classification. Mult-Adds indicate the number of composite multi-add operations per input pixel. Note that some models may contain multiple accurate operations are calculated for the image size reported in the table. Model size for [14] calculated from open-source implementation.

Instance Segmentation



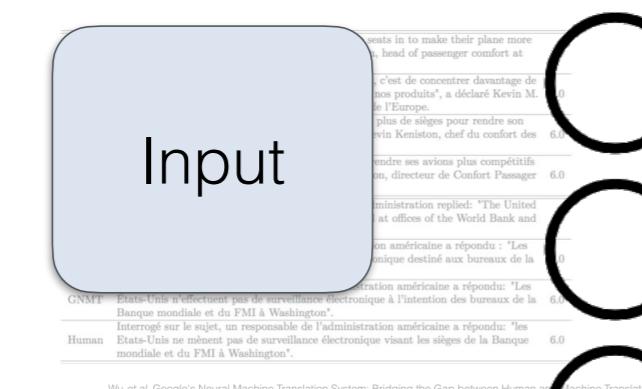
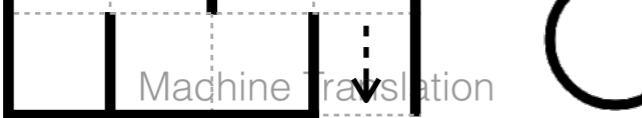
Zoph et al. Learning Transferable Architectures for Scalable Image Recognition

Visual Question Answering

Method	VQA v2 test-dev			VQA v2 test-old				
	All	Yes/no	Numb.	Other	All	Yes/no	Numb.	Other
Prior most common answer in training set [17]	-	-	-	-	25.08	19.08	3.17	3.17
LSTM Language only (blind model) [14]	-	-	-	-	44.26	47.40	31.55	27.37
Deeper LSTM Q score [17] as reported in [14]	-	-	-	-	54.22	73.40	35.18	41.83
MCB [14] as reported in [14]	-	-	-	-	62.27	78.82	38.28	53.36
UPMC-LPBP [1]	-	-	-	-	65.71	82.07	41.06	57.12
Proposed model	-	-	-	-	67.29	83.77	47.97	59.97
HDI-USyd-UNCC	-	-	-	-	68.77	81.80	46.29	58.30
Proposed model	-	-	-	-	68.09	84.50	45.39	59.01
ResNet features T×7, single network	62.07	79.20	39.46	52.62	62.27	79.32	39.77	52.59
Image features from bottom-up attention, adaptive K , single network	65.32	81.82	44.21	56.05	65.67	82.20	43.90	56.26
ResNet features T×7, ensemble	66.34	83.58	43.19	57.40	66.73	83.71	43.77	57.20
Image features from bottom-up attention, adaptive K , ensemble	69.87	86.08	48.09	60.80	70.34	86.64	48.64	61.15

Table 3. Comparison of our best model with competing methods. Excerpt from the official VQA v2 Leaderboard [1].

Teney et al. Tips and Tricks for Visual Question Answering: Learnings from the 2017 Challenge



Wu et al. Google’s Neural Machine Translation System: Bridging the Gap between Human and Machine Translation

Question Answering

System	Dev EM	Dev F1	Test EM	Test F1
Leaderboard (Oct 8th, 2018)				
Human	-	-	82.1	91.2
#1 Ensemble - nfinet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	91.5
#1 Single - nfinet	-	-	83.5	90.3
#2 Single - QANet	-	-	82.5	89.3
Published				
BIDAF+ELMo (Single)	-	85.8	-	-
R.M. Reader (Single)	78.9	86.3	79.5	86.6
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours	80.8	88.5	-	-
BERT _{BASE} (Single)	84.1	90.9	-	-
BERT _{LARGE} (Single)	85.8	91.8	-	-
BERT _{LARGE} (Ensemble)	84.2	91.5	85.1	91.8
BERT _{LARGE} (Sgl.+TriviaQnA)	86.2	92.2	87.4	93.2

Table 2: SQuAD results. The BERT ensemble is 7x systems which use different pre-training checkpoints and fine-tuning seeds.

Devlin et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Word Embeddings



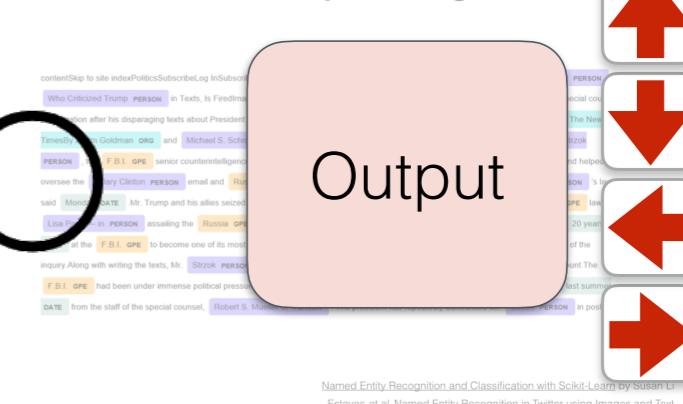
On word embeddings - Part 1 by Ruder

Speech Recognition

Senone set	Model/combination step	WER		WER	
		devset	test	ngram-LM	LSTM-LMs
9k	BLSTM	11.5	8.3	9.2	6.3
27k	BLSTM	11.4	8.0	9.2	6.3
9k-puhpum	BLSTM	11.3	8.0	9.2	5.4
9k-puhpum	BLSTM+ResNet+LACE+CNN-BLSTM	9.6	7.2	7.7	5.4
9k-puhpum	BLSTM+ResNet+LACE	9.7	7.4	7.8	5.4
27k	BLSTM+ResNet+LACE	9.7	7.3	7.8	5.5
	BLSTM+ResNet+LACE	10.0	7.5	8.0	5.8
	Confusion network combination			7.4	5.2
	+ LSTM rescoring			7.3	5.2
	+ ngram rescoring			7.2	5.2
	+ backchannel penalty			7.2	5.1

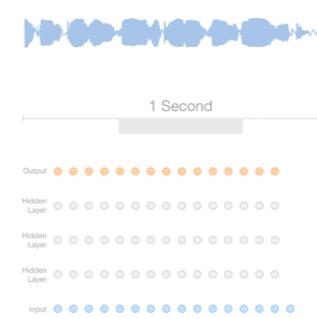
Xiong et al. The Microsoft 2017 Conversational Speech Recognition System

Named Entity Recognition



Named Entity Recognition and Classification with Scikit-Learn by Susan Li
Esteves et al. Named Entity Recognition in Twitter using Images and Text

Speech Synthesis (text-to-speech)



Van Den Oord et al. WaveNet: A Generative Model for Raw Audio

Supervised Learning

Image Classification

Model	image size	# parameters	Mult-Adds	Top 1 Acc. (%)	Top 5 Acc. (%)
Inception V2 [70]	299×299	11.5M	1.64B	74.8	92.2
NASNet-A (5@ 1538)	299×299	6.9M	2.35 B	78.6	94.2
Inception V3 [50]	299×299	11.5M	5.72 B	78.8	94.4
Xception [9]	299×299	11.5M	8.38 B	79.0	94.5
Inception ResNet V2 [58]	299×299	11.5M	13.2 B	80.1	95.1
NASNet-A (7@ 1920)	299×299	11.5M	4.93 B	80.8	95.3
MobileNet v2 [69]	320×320	8.3M	0.27 B	80.9	95.6
PolyNet [69]	333×333	92 M	34.7 B	81.3	95.5
DPN-131 [8]	320×320	79.5M	32.0 B	81.5	95.5
SENet [25]	320×320	145.8M	42.3 B	82.7	95.6
NASNet-A (6@ 4032)	333×333	88.9 M	23.8 B	82.7	95.2

Table 1. Performance of architecture search and other published state-of-the-art models on ImageNet classification. Mult-Adds indicate the number of composite multi-addition operations per input pixel. Note that some models may contain multiple multi-addition operations are calculated for the image size reported in the table. Model size for [14] is calculated from open-source implementation.

Instance Segmentation



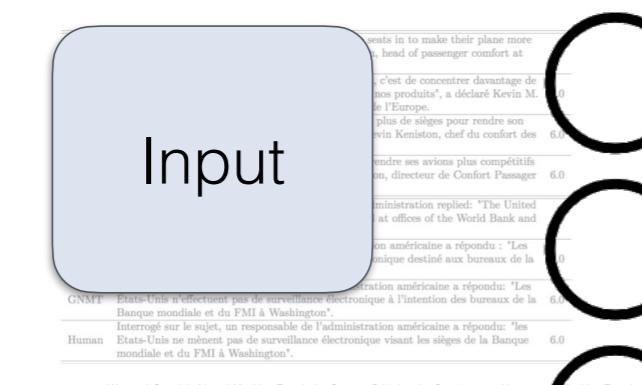
Zoph et al. Learning Transferable Architectures for Scalable Image Recognition

Visual Question Answering

Method	VQA v2 test-dev		
	All	Yes/no	Num. Other
Prior most common answer in training set [1]	-	-	25.00 39.00 37.00
LSTM Language only (blind model) [1]	-	-	44.26 47.40 31.55 27.37
Deeper LSTM Q score [17] as reported in [1]	-	-	54.22 73.40 35.18 41.83
MCB [10] as reported in [1]	-	-	62.27 78.82 38.28 53.36
UPMC-LPBP [1]	-	-	65.71 82.07 41.06 57.12
Proposed model	-	-	67.29 83.77 40.97 59.97
HDI-USyd-UNCC	-	-	68.09 84.50 45.39 59.01
Proposed model	-	-	68.09 84.50 45.39 59.01
ResNet features T×7, single network	62.07	79.20	39.46 52.62
Image features from bottom-up attention, adaptive K , single network	65.32	81.82	44.21 56.05
ResNet features T×7, ensemble	66.34	83.58	43.19 57.40
Image features from bottom-up attention, adaptive K , ensemble	69.87	86.08	48.09 60.80

Table 3. Comparison of our best model with competing methods. Excerpt from the official VQA v2 Leaderboard [1].

Teney et al. Tips and Tricks for Visual Question Answering: Learnings from the 2017 Challenge



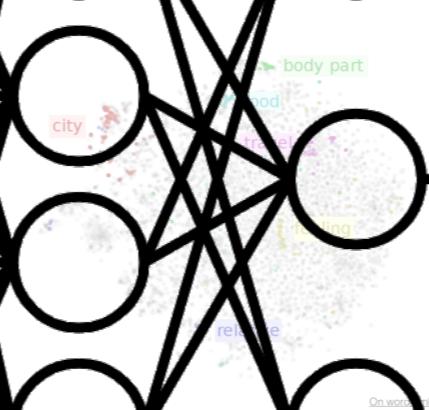
Question Answering

System	Dev EM	Dev F1	Test EM	Test F1
Leaderboard (Oct 8th, 2018)				
Human	-	-	82.1	91.2
#1 Ensemble - nlnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	91.5
#1 Single - nlnet	-	-	83.5	90.3
#2 Single - QANet	-	-	82.5	89.3
Published				
BIDAF+ELMo (Single)	-	85.8	-	-
R.M. Reader (Single)	78.9	86.3	79.5	86.6
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours	80.8	88.5	-	-
BERT _{BASE} (Single)	84.1	90.9	-	-
BERT _{LARGE} (Single)	85.8	91.8	-	-
BERT _{LARGE} (Ensemble)	84.2	91.5	85.1	91.8
BERT _{LARGE} (Sgl.+TriviaQnA)	86.2	92.2	87.4	93.2

Table 2: SQuAD results. The BERT ensemble is 7x systems which use different pre-training checkpoints and fine-tuning seeds.

Devlin et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Word Embeddings

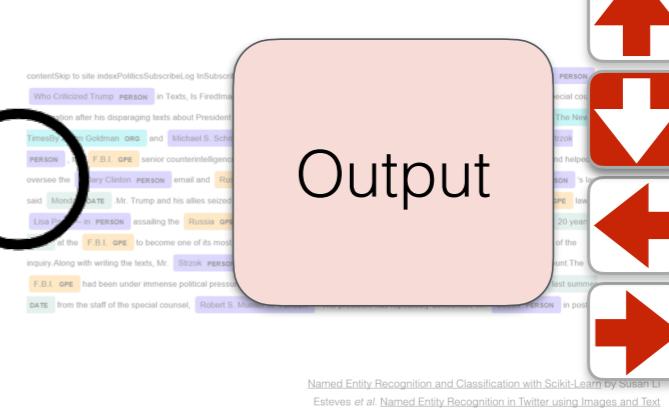


Speech Recognition

Senone set	Model/combination step	WER		WER	
		devset	test ngram-LM	devset	test LSTM-LMs
9k	BLSTM	11.5	8.3	9.2	6.3
27k	BLSTM	11.4	8.0	9.2	6.3
9k-puhpum	BLSTM	11.3	8.0	9.2	6.3
9k-puhpum	BLSTM+ResNet+LACE+CNN-BLSTM	9.6	7.2	7.7	5.4
9k-puhpum	BLSTM+ResNet+LACE	9.7	7.4	7.8	5.4
27k	BLSTM+ResNet+LACE	9.7	7.3	7.8	5.5
	BLSTM+ResNet+LACE	10.0	7.5	8.0	5.8
	Confusion network combination			7.4	5.2
	+ LSTM rescoring			7.3	5.2
	+ ngram rescoring			7.2	5.2
	+ backchannel penalty			7.2	5.1

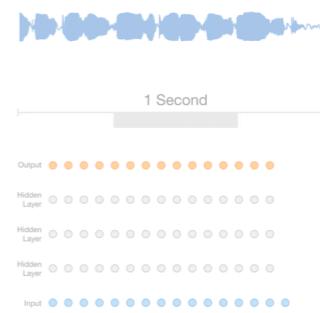
Word Error Rate

Named Entity Recognition



Named Entity Recognition and Classification with Scikit-Learn by Susan Li
Esteves et al. Named Entity Recognition in Twitter using Images and Text

Speech Synthesis (text-to-speech)



Van Den Oord et al. WaveNet: A Generative Model for Raw Audio

Xiong et al. The Microsoft 2017 Conversational Speech Recognition System

Supervised Learning

Image Classification

Model	image size	# parameters	Mult-Adds	Top 1 Acc. (%)	Top 5 Acc. (%)
Inception V2 [70]	299×299	11.5M	1.64B	74.8	92.2
NASNet-A (5@ 1538)	299×299	10.9 M ^b	2.35 B	78.6	94.2
Inception V3 [50]	299×299	23.8M ^b	5.72 B	78.8	94.4
Xception [9]	299×299	22.8M ^b	8.38 B	79.0	94.5
Inception ResNet V2 [58]	299×299	55.8M ^b	13.2 B	80.1	95.1
NASNet-A (7@ 1920)	299×299	22.6M ^b	4.93 B	80.8	95.3
ResNet-50 [14]	224×224	23.6M ^b	2.02 B	80.9	95.6
PolyNet [69]	333×333	1.5 M ^b	34.7 B	81.3	95.5
DPN-131 [8]	320×320	1.5 M ^b	32.0 B	81.5	95.5
SENet [25]	320×320	1.5 M ^b	42.3 B	82.7	95.6
NASNet-A (6@ 4032)	333×333	1.5 M ^b	23.8 B	82.7	95.2

Table 1. Performance of architecture search and other selected state-of-the-art models on ImageNet classification. Mult-Adds indicate the number of composite multi-add operations per input pixel. Note that some models may contain multiple accumulate operations are calculated for the image size reported in the table. Model size for [14] is calculated from open-source implementation.

Zoph et al. Learning Transferable Architectures for Scalable Image Recognition

Instance Segmentation



Zhou et al. Mask R-CNN: Real-Time Multi-Object Instance Segmentation

Visual Question Answering

Method	VQA v2 test-dev		
	All	Yes/no	Num. Other
Prior (most common answer in training set) [1]	-	-	25.00 39.00 37.00
LSTM Language only (blind model) [1]	-	-	44.26 47.40 31.55 27.37
Deeper LSTM Q & vnn. [17] as reported in [1]	-	-	54.22 73.40 35.18 41.83
MCB [1] as reported in [1]	-	-	62.27 78.82 38.28 53.36
UPMC-LPBP [1]	-	-	65.71 82.07 41.06 57.12
Proposed model	-	-	67.29 83.77 40.97 59.97
HCU-USyd-UNCC	-	-	68.09 84.50 45.39 59.01
Proposed model	-	-	68.09 84.50 45.39 59.01
ResNet features T×7, single network	62.07	79.20	39.46 52.62
Image features from bottom-up attention, adaptive K , single network	65.32	81.82	44.21 56.05
ResNet features T×7, ensemble	66.34	83.58	43.19 57.40
Image features from bottom-up attention, adaptive K , ensemble	69.87	86.08	48.00 68.64 61.15

Table 3. Comparison of our best model with competing methods. Excerpt from the official VQA v2 Leaderboard [1].

Teney et al. Tips and Tricks for Visual Question Answering: Learnings from the 2017 Challenge

Zoph et al. Learning Transferable Architectures for Scalable Image Recognition

He et al. Mask R-CNN

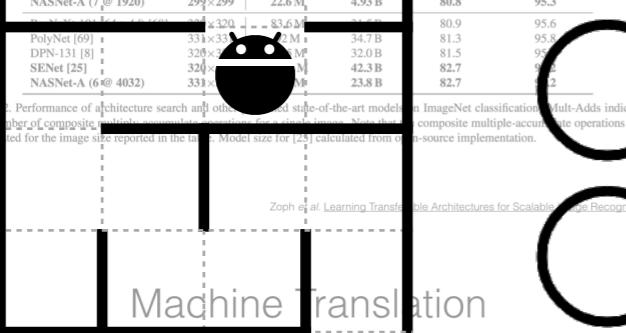
He et al. Mask

Supervised Learning

Image Classification

Model	image size	# parameters	Mult-Adds	Top 1 Acc. (%)	Top 5 Acc. (%)
Inception V2 [70]	299×299	11.5M	1.64B	74.8	92.2
NASNet-A (5@ 1538)	299×299	10.9 M ^b	2.35 B	78.6	94.2
Inception V3 [50]	299×299	23.8M ^b	5.72 B	78.8	94.4
Xception [9]	299×299	22.8M ^b	8.38 B	79.0	94.5
Inception ResNet V2 [58]	299×299	55.8M ^b	13.2 B	80.1	95.1
NASNet-A (7@ 1920)	299×299	22.6M ^b	4.93 B	80.8	95.3
ResNet-50 [14]	224×224	23.6M ^b	2.02 B	80.9	95.6
PolyNet [69]	333×333	1.5 M ^b	34.7 B	81.3	95.5
DPN-131 [8]	320×320	1.5 M ^b	32.0 B	81.5	95.5
SENet [25]	320×320	1.5 M ^b	42.3 B	82.7	95.6
NASNet-A (6@ 4032)	331×331	1.5 M ^b	23.8 B	82.7	95.2

Table 1. Performance of architecture search and other selected state-of-the-art models on ImageNet classification. Mult-Adds indicate the number of composite multi-add operations per input pixel. Note that some models may contain multiple accumulate operations are calculated for the image size reported in the table. Model size for [14] is calculated from open-source implementation.



Instance Segmentation

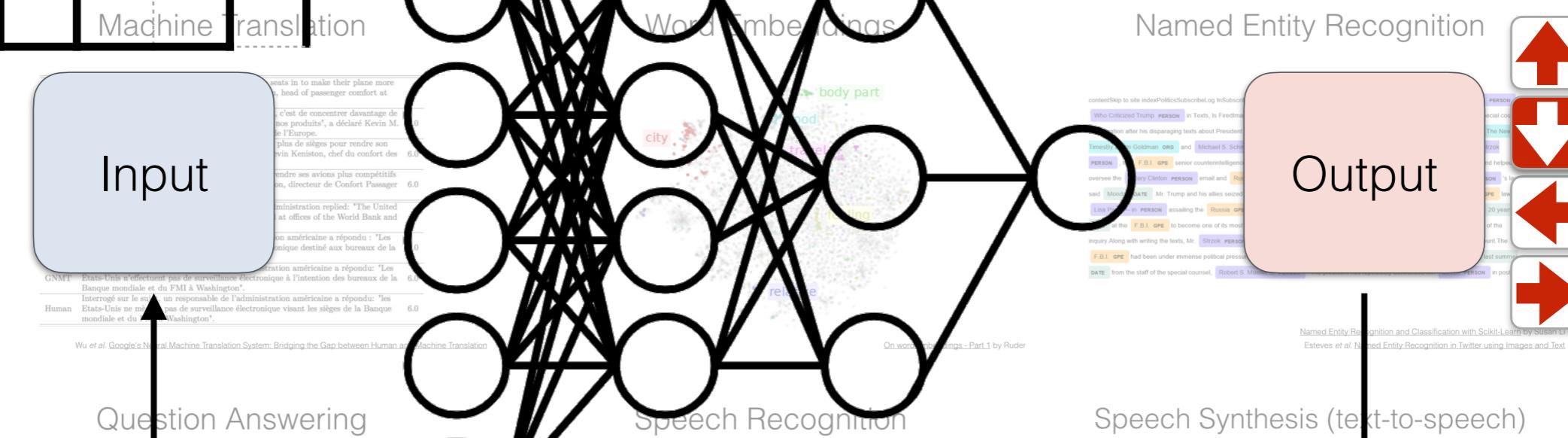


Visual Question Answering

Method	VQA v2 test-dev			VQA v2 test-old				
	All	Yes/no	Numth.	Other	All	Yes/no	Numth.	Other
Prior most common answer in training set [1]	—	—	—	—	25.00	19.00	31.55	3.17
LSTM Language only (blind model) [1]	—	—	—	—	44.26	47.01	31.55	27.37
Deeper LSTM Q score [17] as reported in [1]	—	—	—	—	54.22	73.40	35.18	41.83
MCB [11] as reported in [1]	—	—	—	—	62.27	78.82	38.28	53.36
UPMC-LPBP [1]	—	—	—	—	65.71	82.07	41.06	57.12
UNICUS [18]	—	—	—	—	67.29	82.07	41.06	57.17
HDI-USyd-UNCC [19]	—	—	—	—	68.77	81.89	46.29	58.30
Proposed model	—	—	—	—	68.09	84.50	45.39	59.01

Table 3. Comparison of our best model with competing methods. Excerpt from the official VQA v2 Leaderboard [1].

Teney et al. Tips and Tricks for Visual Question Answering: Learnings from the 2017 Challenge



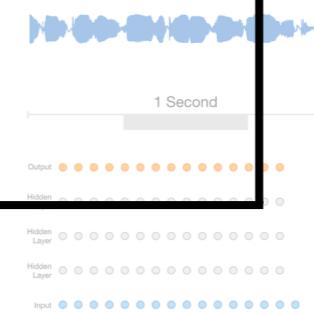
System	Dev EM	Dev F1	Test EM	Test F1
Leaderboard (Oct 8th, 2018)				
Human	-	-	82.9	91.2
#1 Ensemble - nInet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	91.5
#1 Single - nInet	-	-	83.5	90.3
#2 Single - QANet	-	-	82.5	89.3
Published				
BIDAF+ELMo (Single)	-	85.8	-	-
R.M. Reader (Single)	78.9	86.3	79.5	86.6
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours	80.8	88.5	-	-
BIDAF+ELMo (Ensemble)	84.1	90.0	-	-
BERT _{BASE} (Single)	85.8	91.8	-	-
BERT _{BASE} (Ensemble)	85.8	91.8	-	-
BERT _{BASE} (Sgl+TriviaQA)	84.2	91.5	85.1	91.8
BERT _{BASE} (Ens+TriviaQA)	86.2	92.2	87.4	93.2

Table 2: SQuAD results. The BERT ensemble is 7x systems which use different pre-training checkpoints and fine-tuning seeds.

Devlin et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

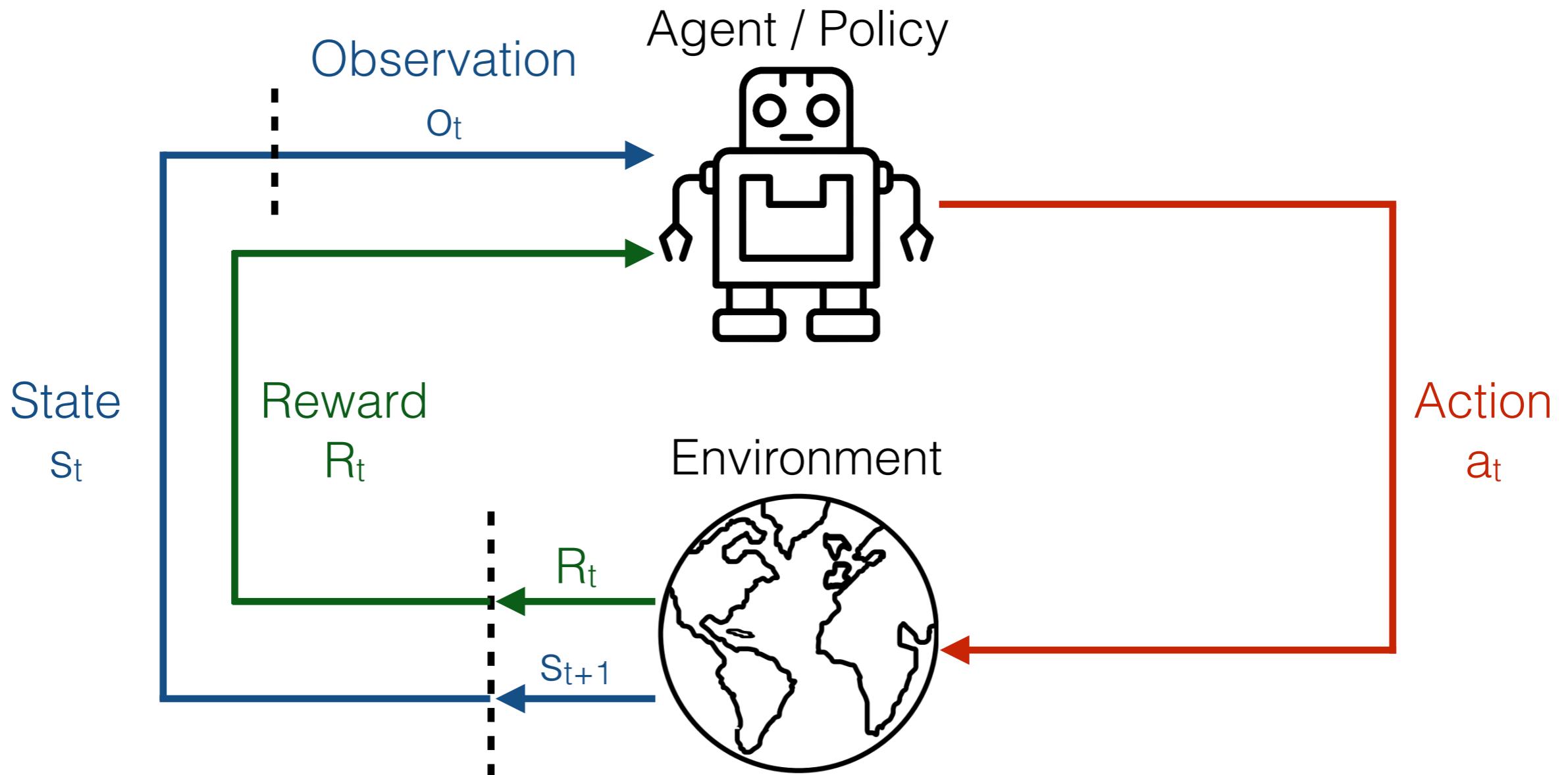
Senone set	Model/combination step	WER		WER	
		devset	test	devset	test
9k	BLSTM	11.5	8.3	9.2	6.3
27k	BLSTM	11.4	8.0	9.2	6.3
9k-puhpum	BLSTM	11.3	8.0	9.2	6.3
9k-puhpum	BLSTM+ResNet+LACE+CNN-BLSTM	9.6	7.2	7.7	5.4
9k-puhpum	BLSTM+ResNet+LACE	9.7	7.4	7.8	5.4
27k	BLSTM+ResNet+LACE	9.7	7.3	7.8	5.5
27k	BLSTM+ResNet+LACE Configurable combination	10.0	7.5	8.0	5.8
		7.4	5.2		
	+ LSTM rescoring			7.3	5.2
	+ ngram rescoring			7.2	5.2
	+ backchannel penalty			7.2	5.1

Xiong et al. The Microsoft 2017 Conversational Speech Recognition System



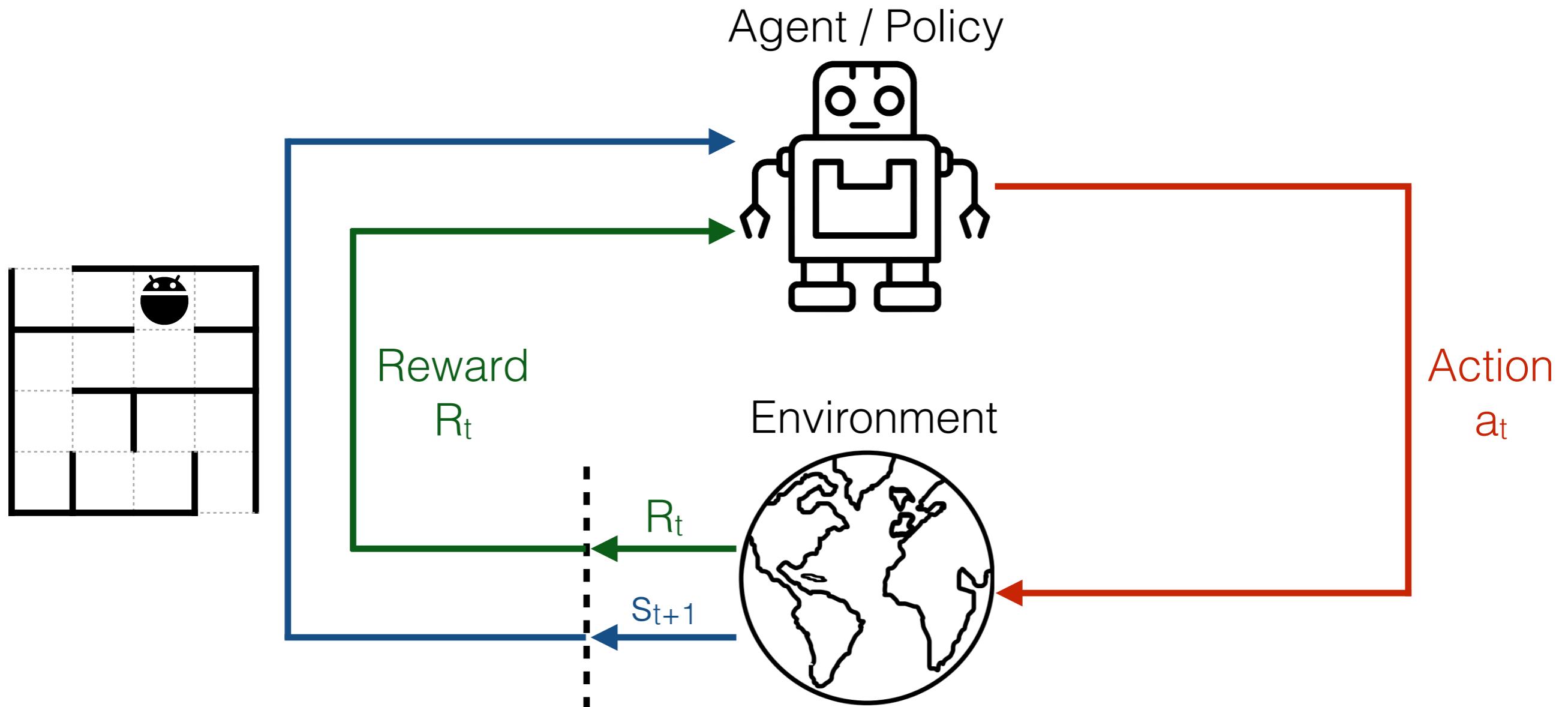
Van Den Oord et al. WaveNet: A Generative Model for Raw Audio

Robot Learning via Reinforcement Learning



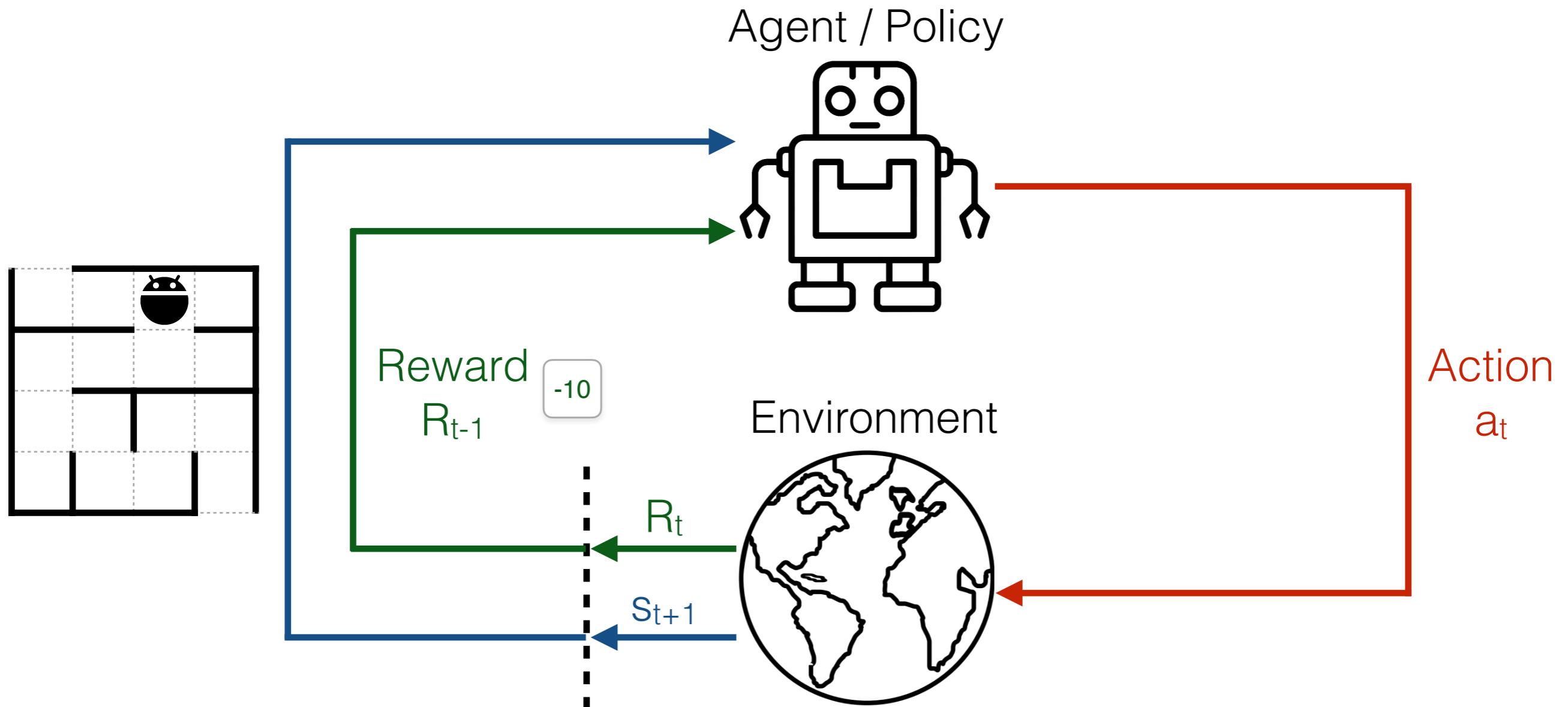
Goal: maximize $\sum_{t=0}^{t=H} \gamma^t R_t(s_t, a_t)$

Robot Learning via Reinforcement Learning



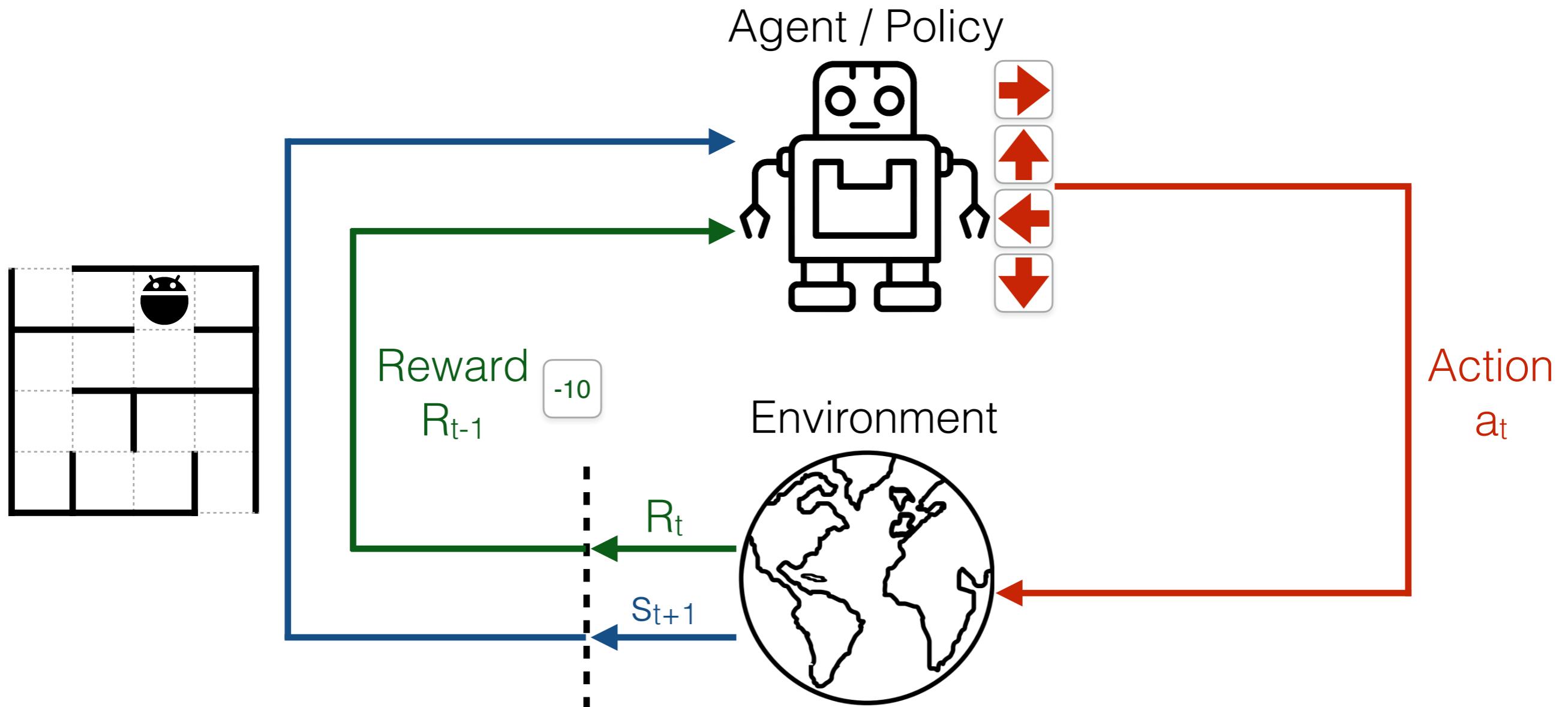
Goal: maximize $\sum_{t=0}^{t=H} \gamma^t R_t(s_t, a_t)$

Robot Learning via Reinforcement Learning



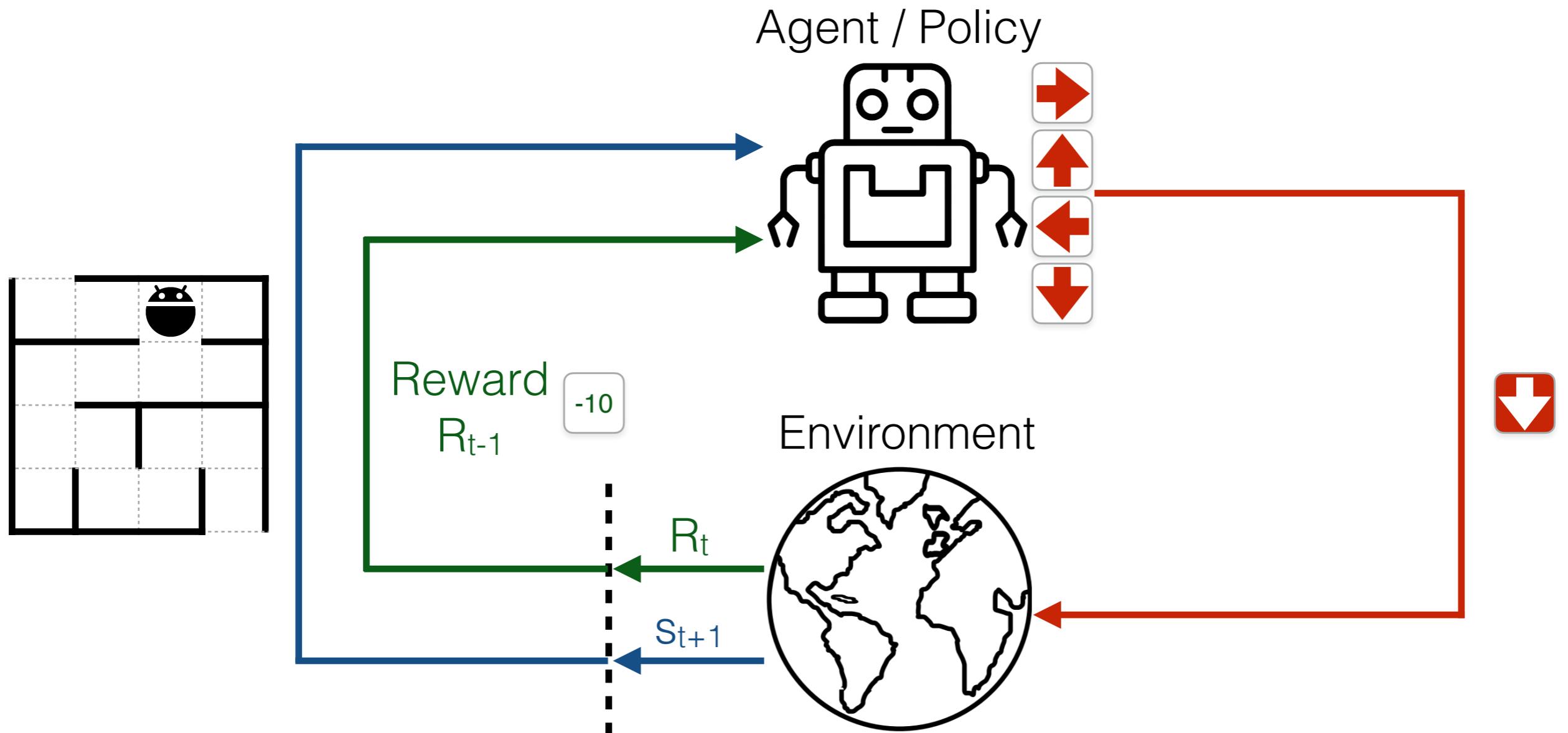
Goal: maximize $\sum_{t=0}^{t=H} \gamma^t R_t(s_t, a_t)$

Robot Learning via Reinforcement Learning



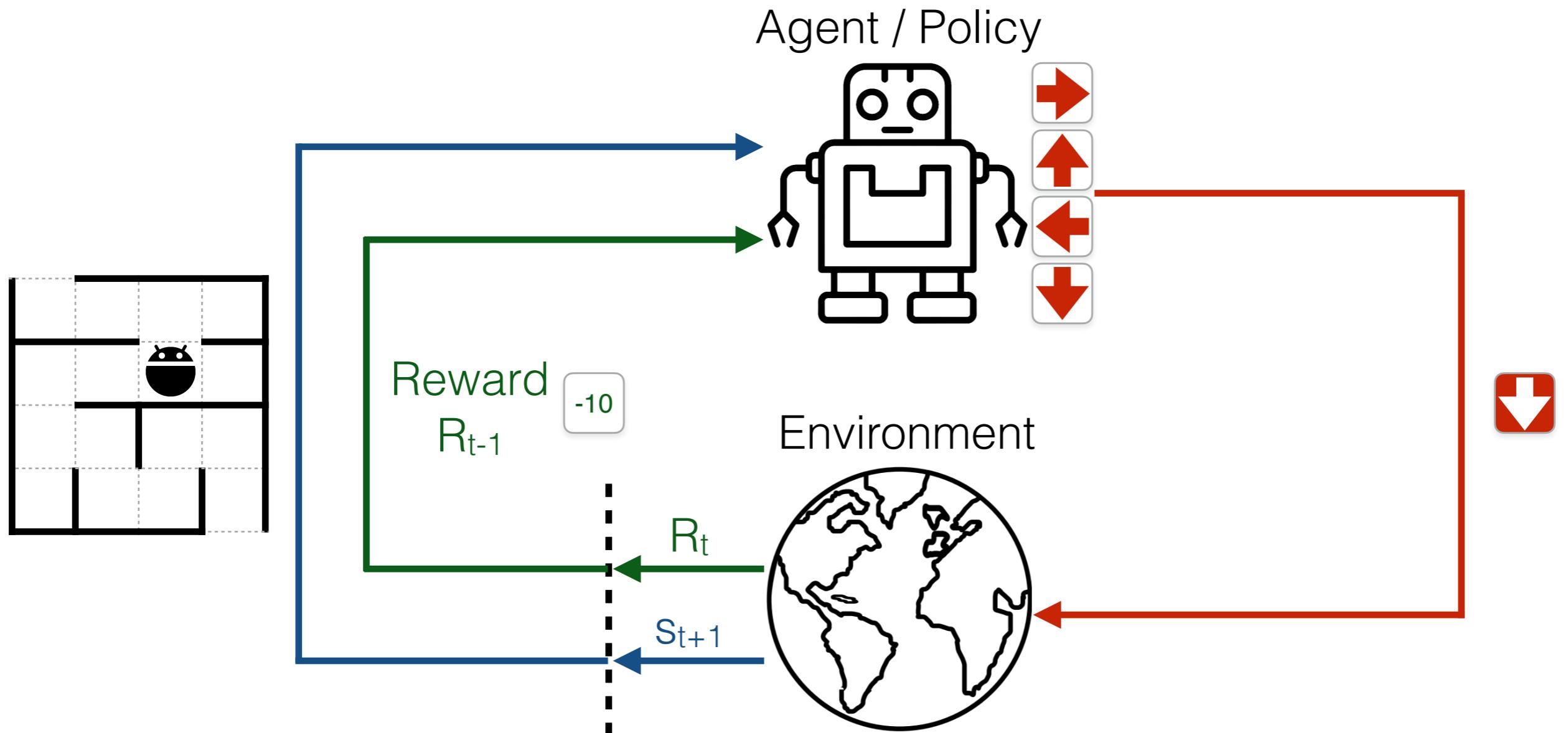
Goal: maximize $\sum_{t=0}^{t=H} \gamma^t R_t(s_t, a_t)$

Robot Learning via Reinforcement Learning



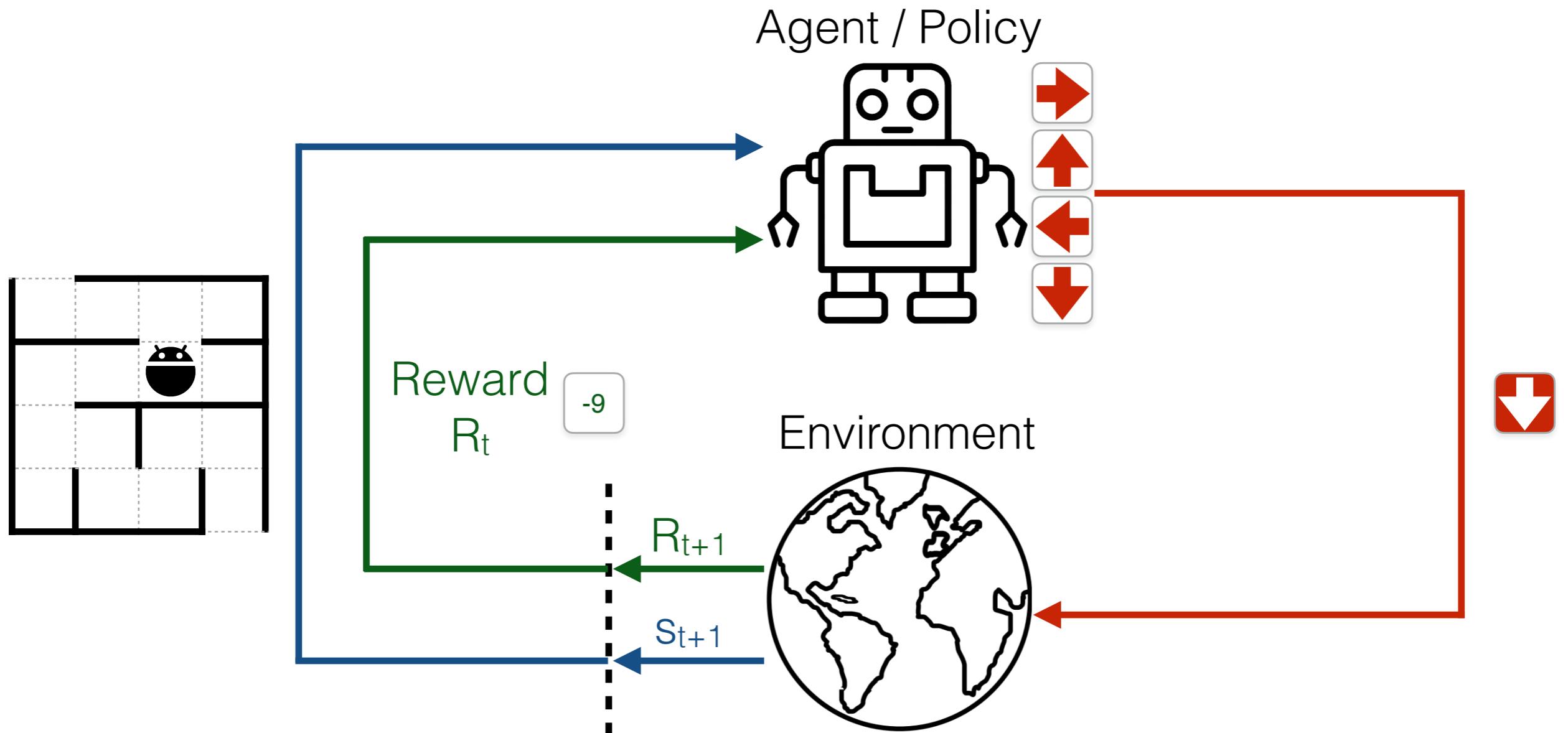
Goal: maximize $\sum_{t=0}^{t=H} \gamma^t R_t(s_t, a_t)$

Robot Learning via Reinforcement Learning



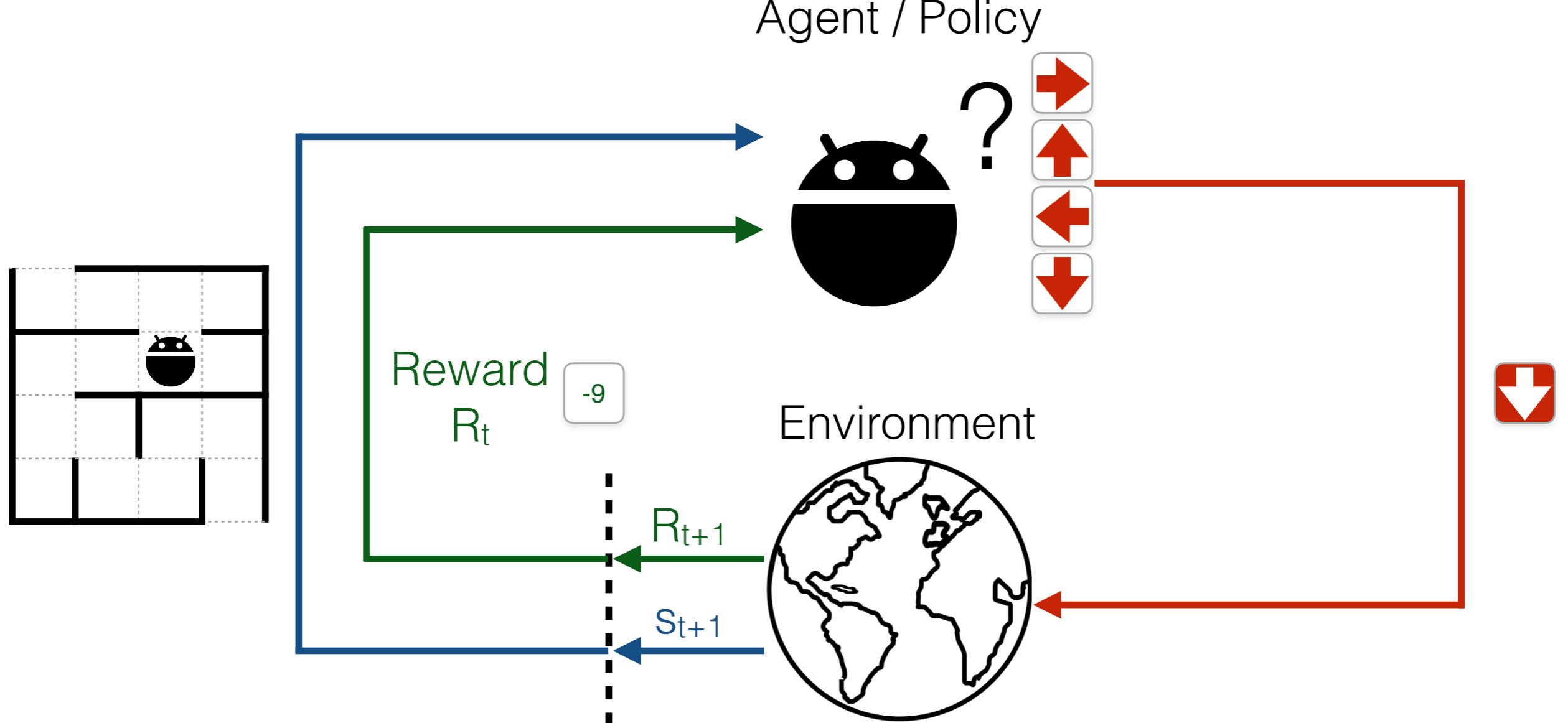
Goal: maximize $\sum_{t=0}^{t=H} \gamma^t R_t(s_t, a_t)$

Robot Learning via Reinforcement Learning



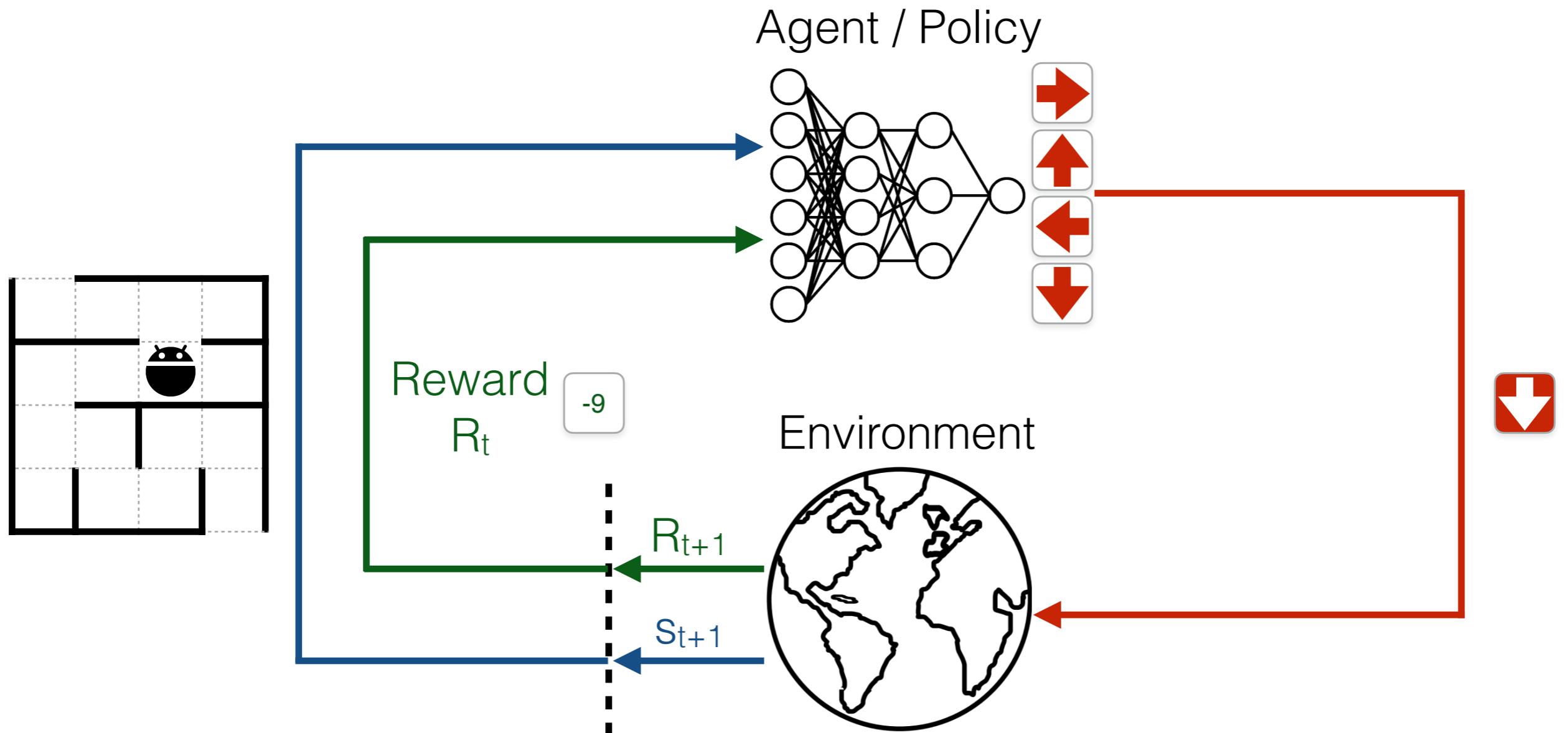
Goal: maximize $\sum_{t=0}^{t=H} \gamma^t R_t(s_t, a_t)$

Robot Learning via Reinforcement Learning



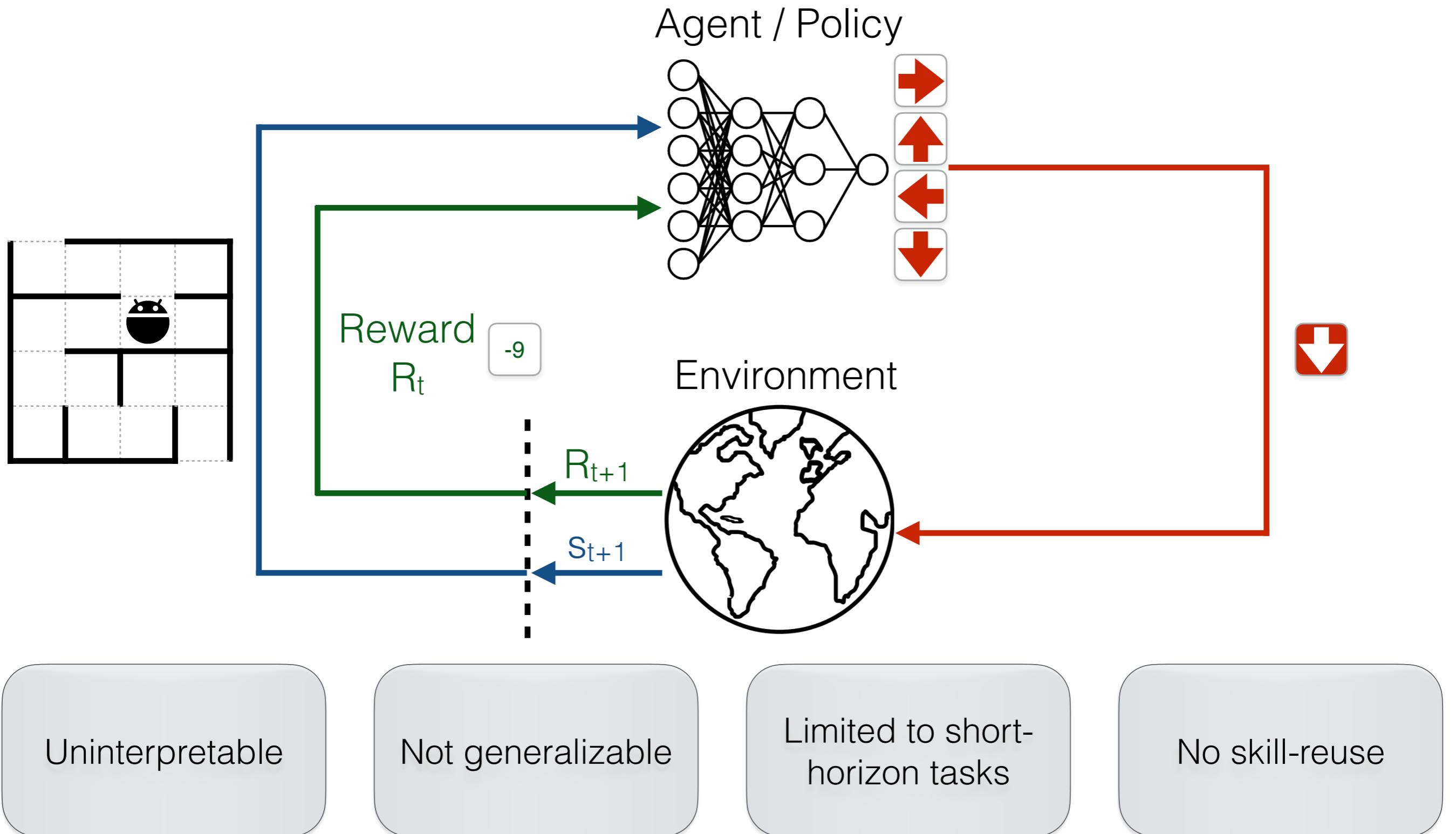
Goal: maximize $\sum_{t=0}^{t=H} \gamma^t R_t(s_t, a_t)$

Robot Learning via Deep Reinforcement Learning

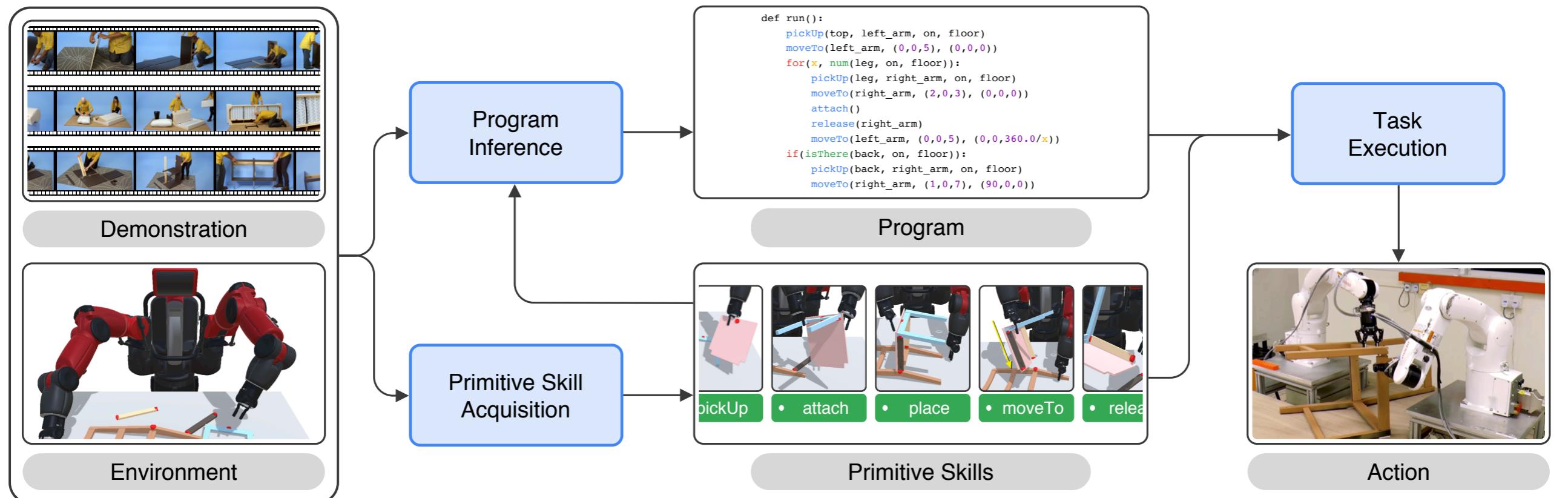


Goal: maximize $\sum_{t=0}^{t=H} \gamma^t R_t(s_t, a_t)$

Robot Learning via Deep Reinforcement Learning



Program-Guided Framework for Interpreting and Acquiring Complex Skills



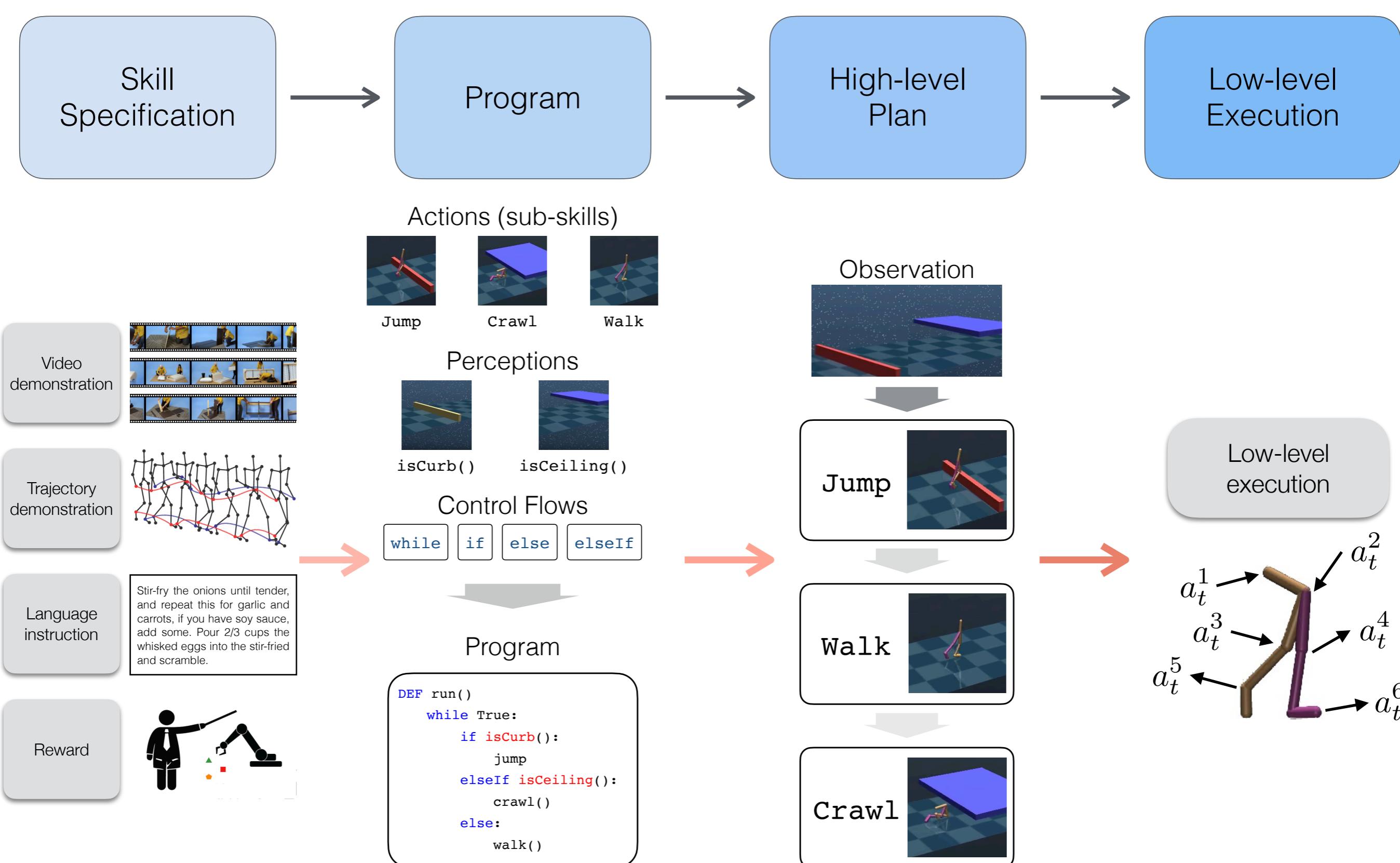
Interpretable

Programmatic / Generalizable

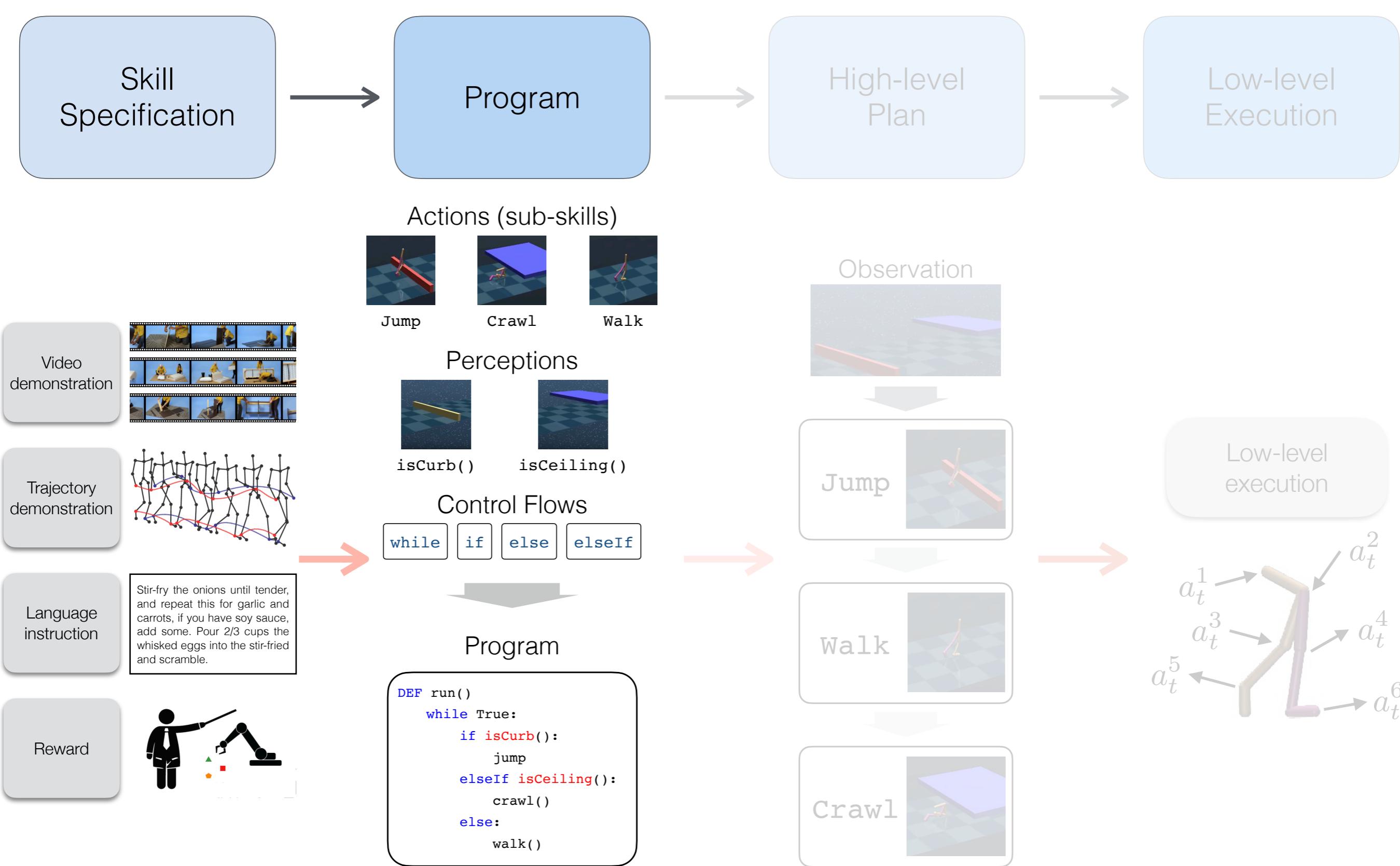
Hierarchical

Modular

Program-Guided Framework for Interpreting and Acquiring Complex Skills

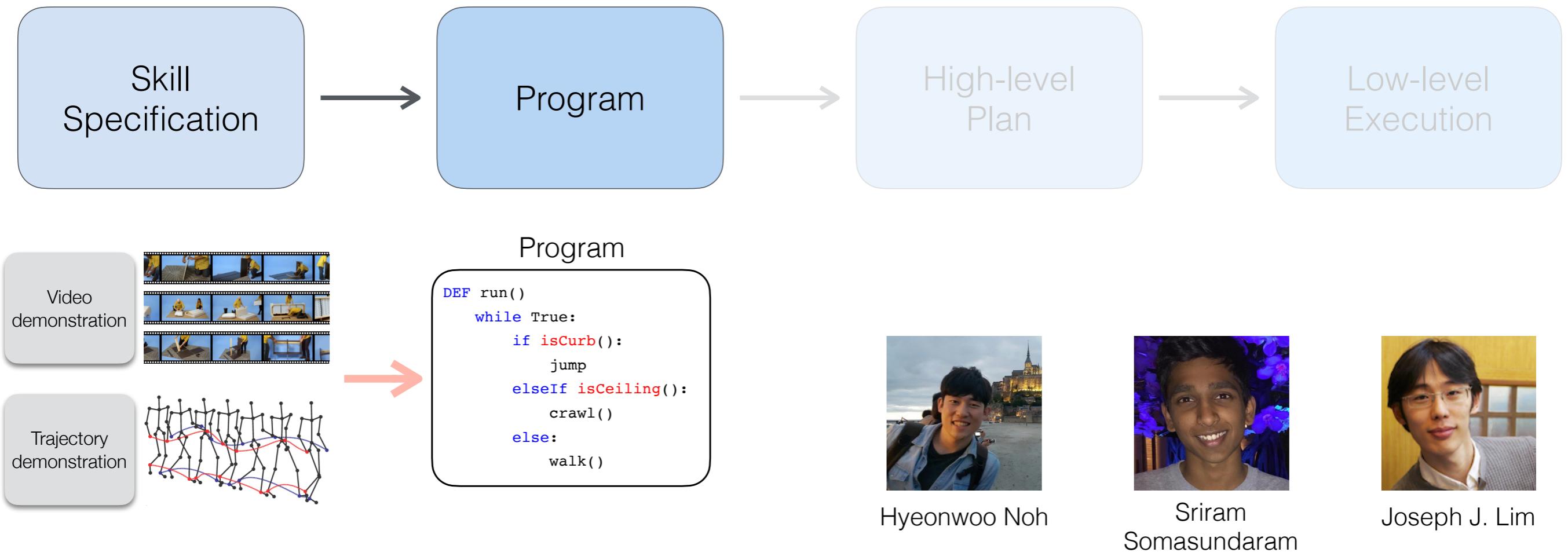


Program Inference



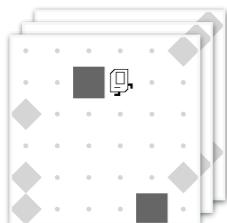
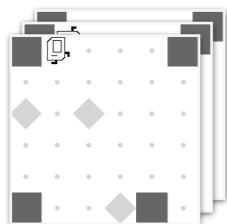
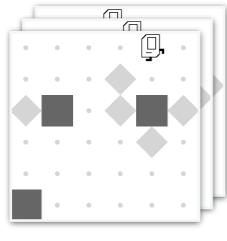
Neural Program Synthesis from Diverse Demonstration Videos

ICML 2018

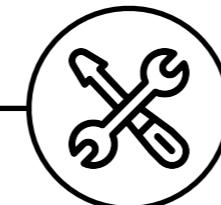


Imitation Learning

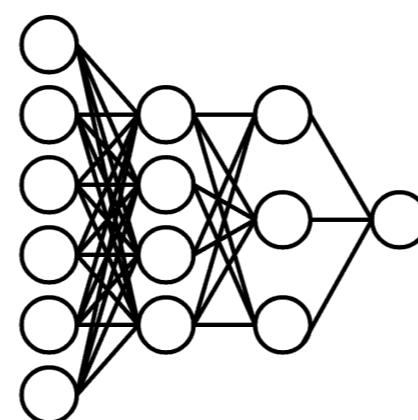
Demonstrations



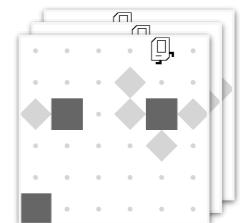
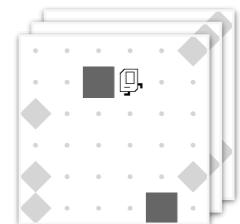
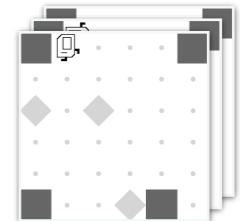
Imitate



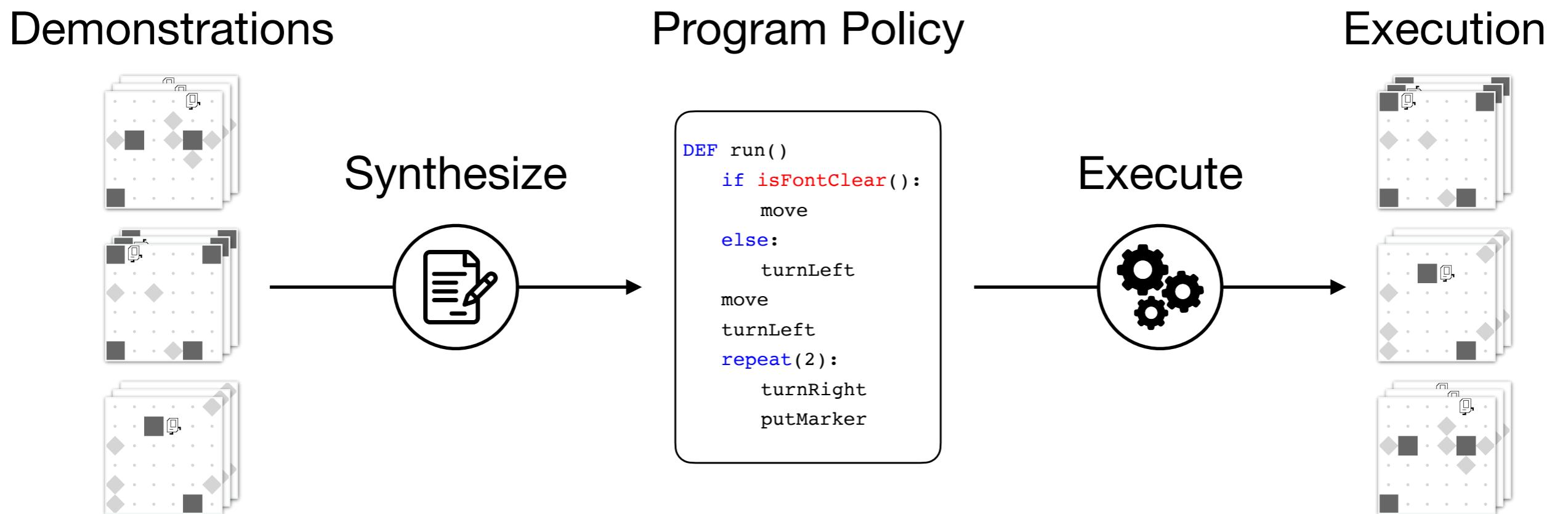
Neural Network
Policy



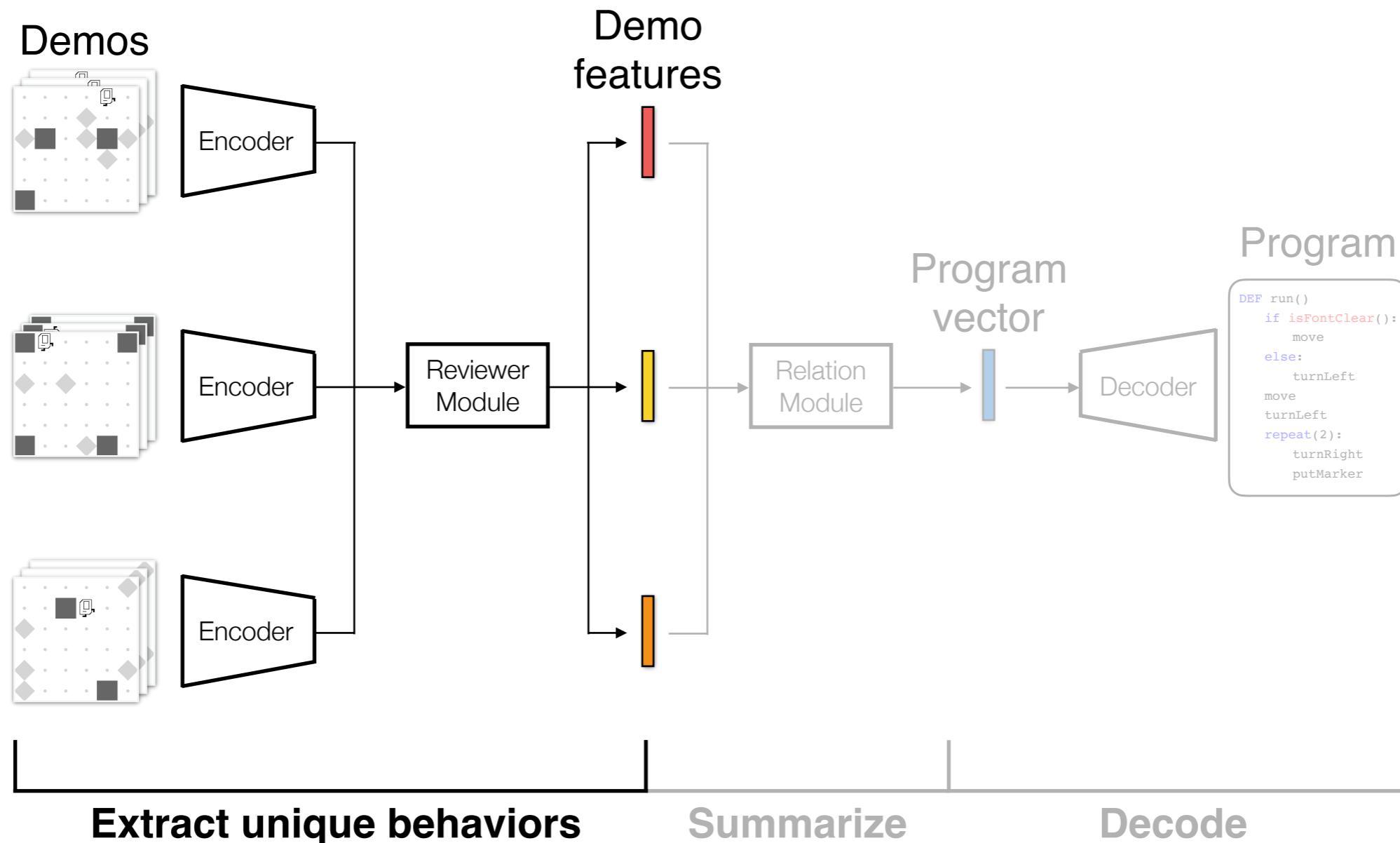
Execution



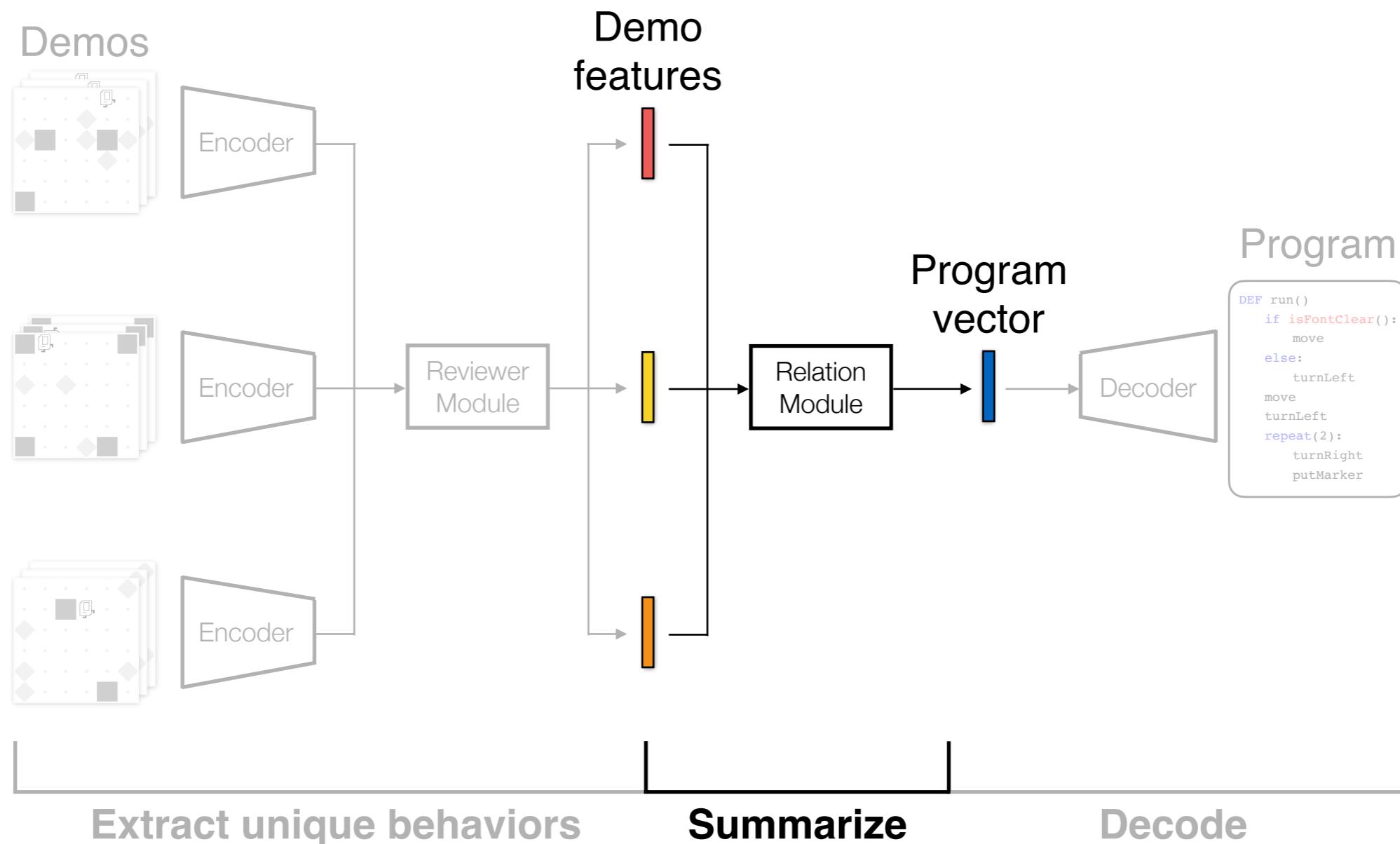
Imitation Learning by Synthesizing Programs



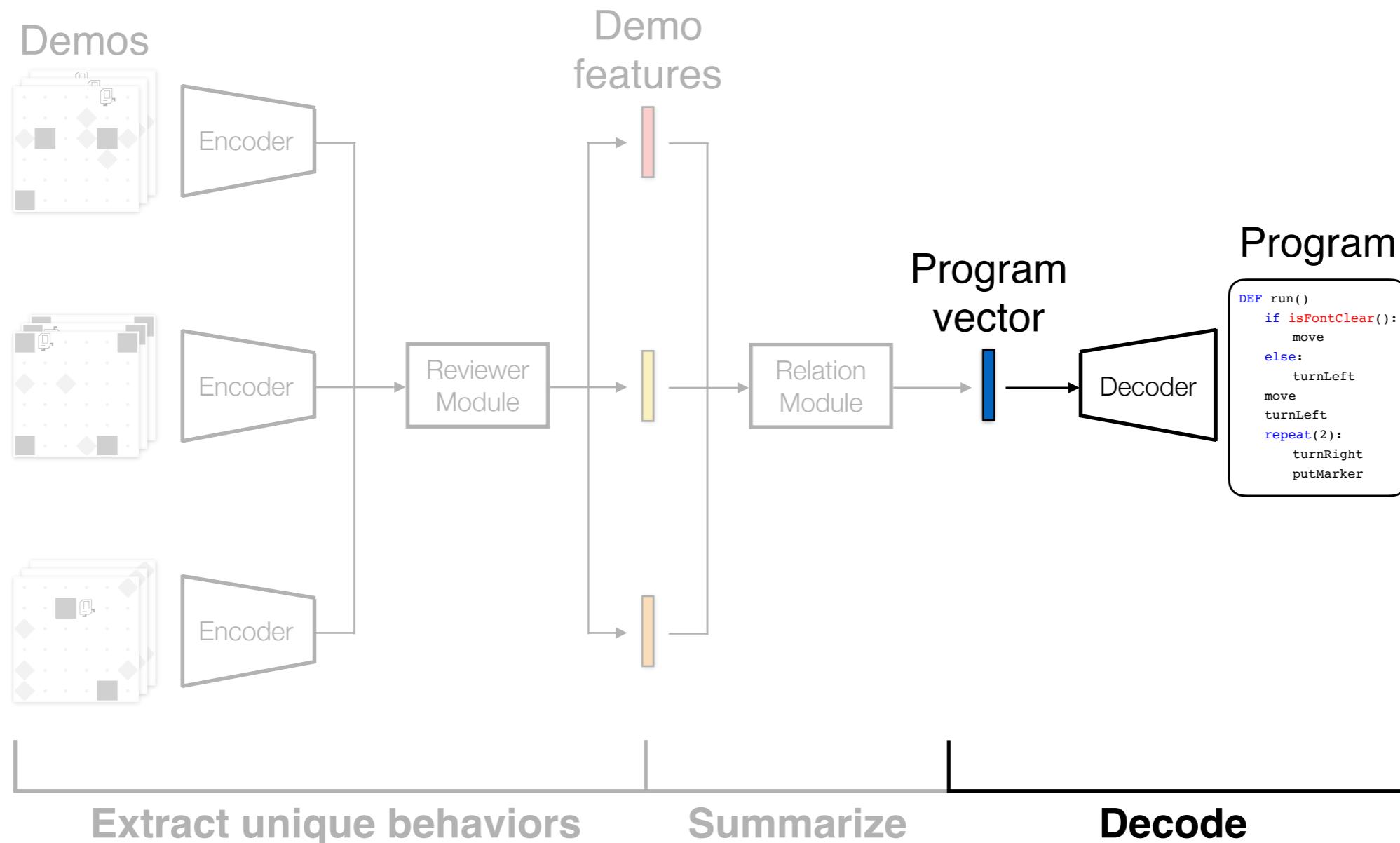
Model Overview



Model Overview



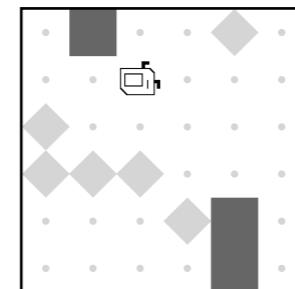
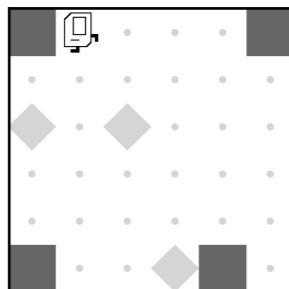
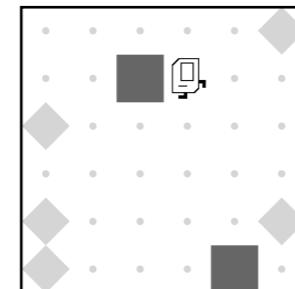
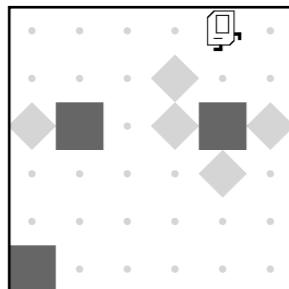
Model Overview



Environments

Karel

```
DEF run()
    if isFontClear():
        move
    else:
        turnLeft
    move
    turnLeft
    repeat(2):
        turnRight
        putMarker
```



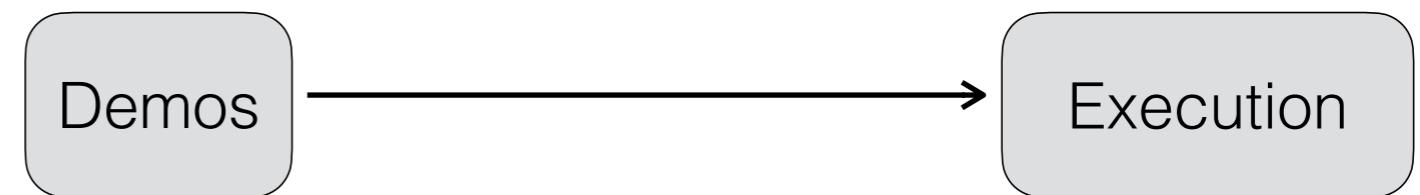
ViZDoom

```
DEF run()
    while isFontClear(HellKnight):
        attack
        moveForward
        if isThere(Demon):
            moveRight
        else:
            moveLeft
            moveBackward
```

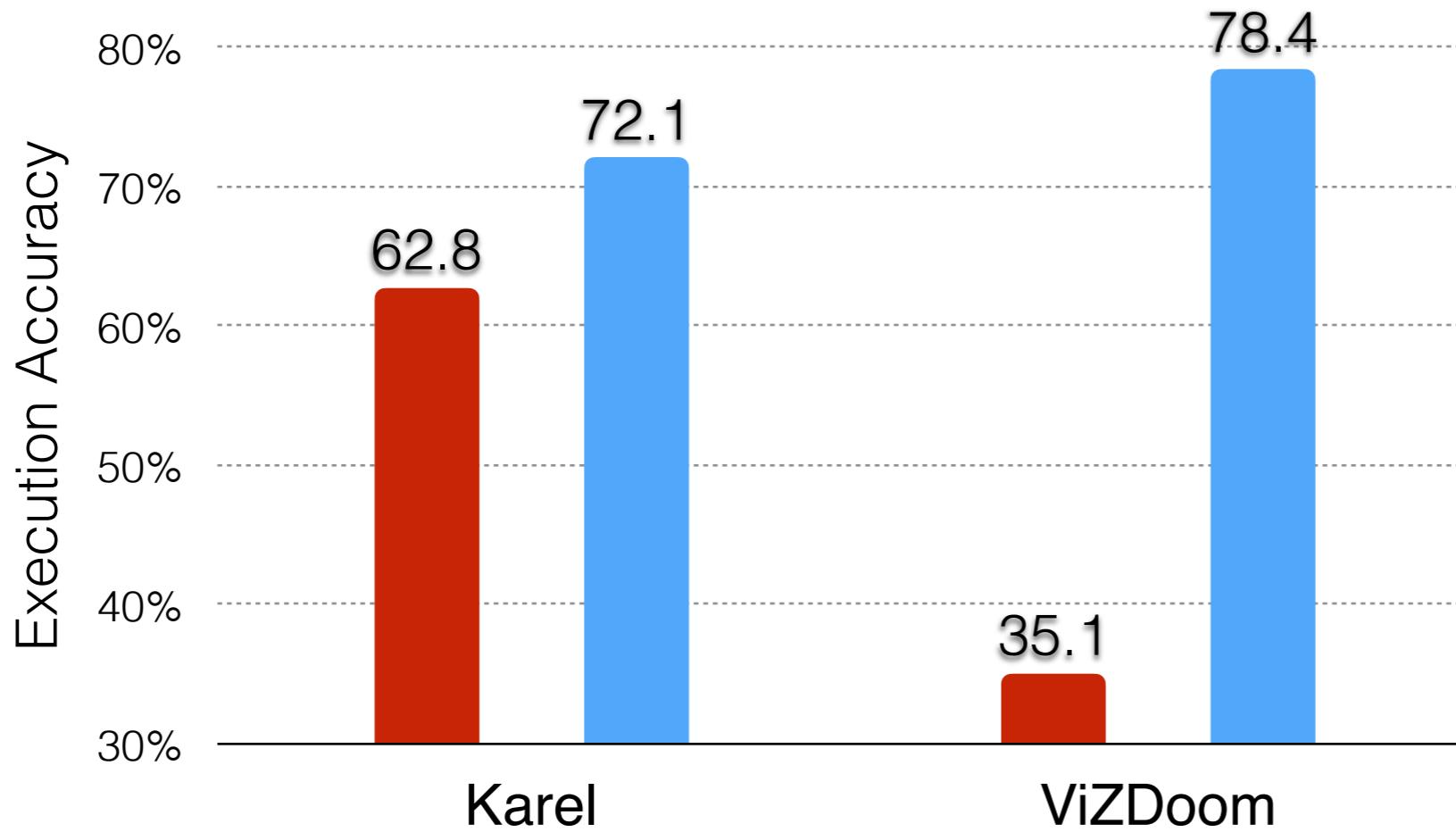
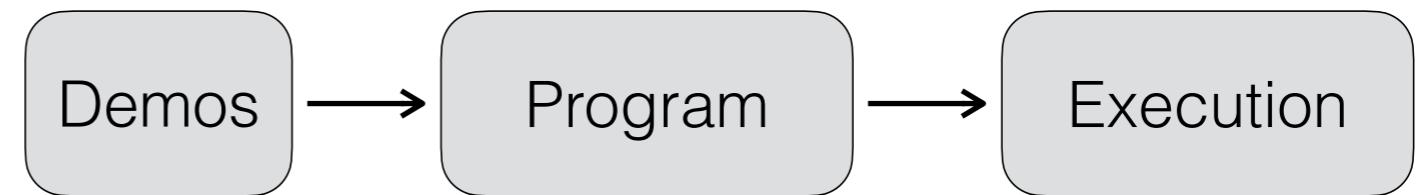


Quantitative Results

Neural Network Policy



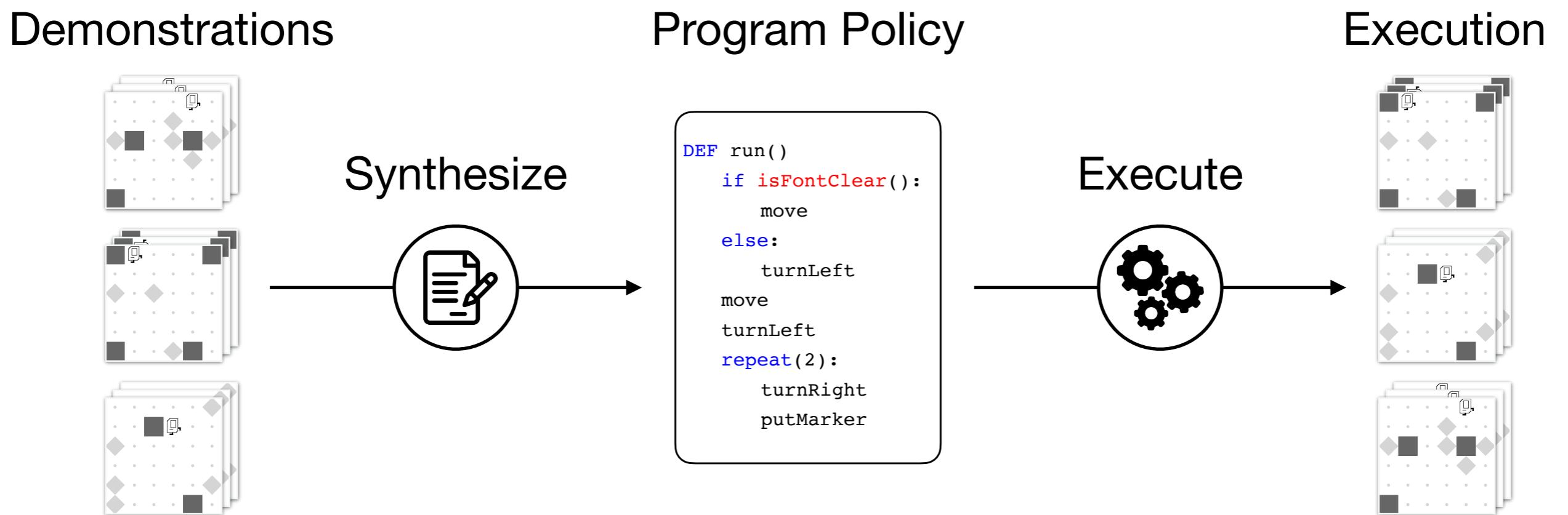
Program Policy



```
DEF run()
    if isFontClear():
        move
    else:
        turnLeft
        move
        turnLeft
        repeat(2):
            turnRight
            putMarker
```

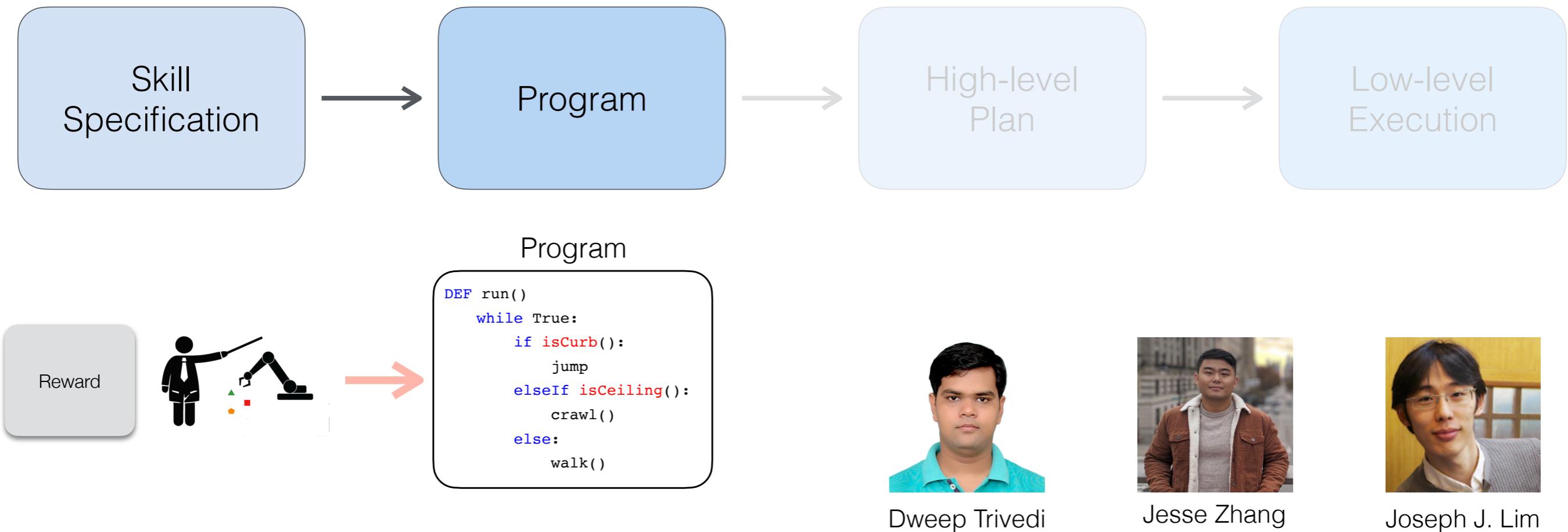
Takeaway

- Synthesize programs to imitate demonstrations

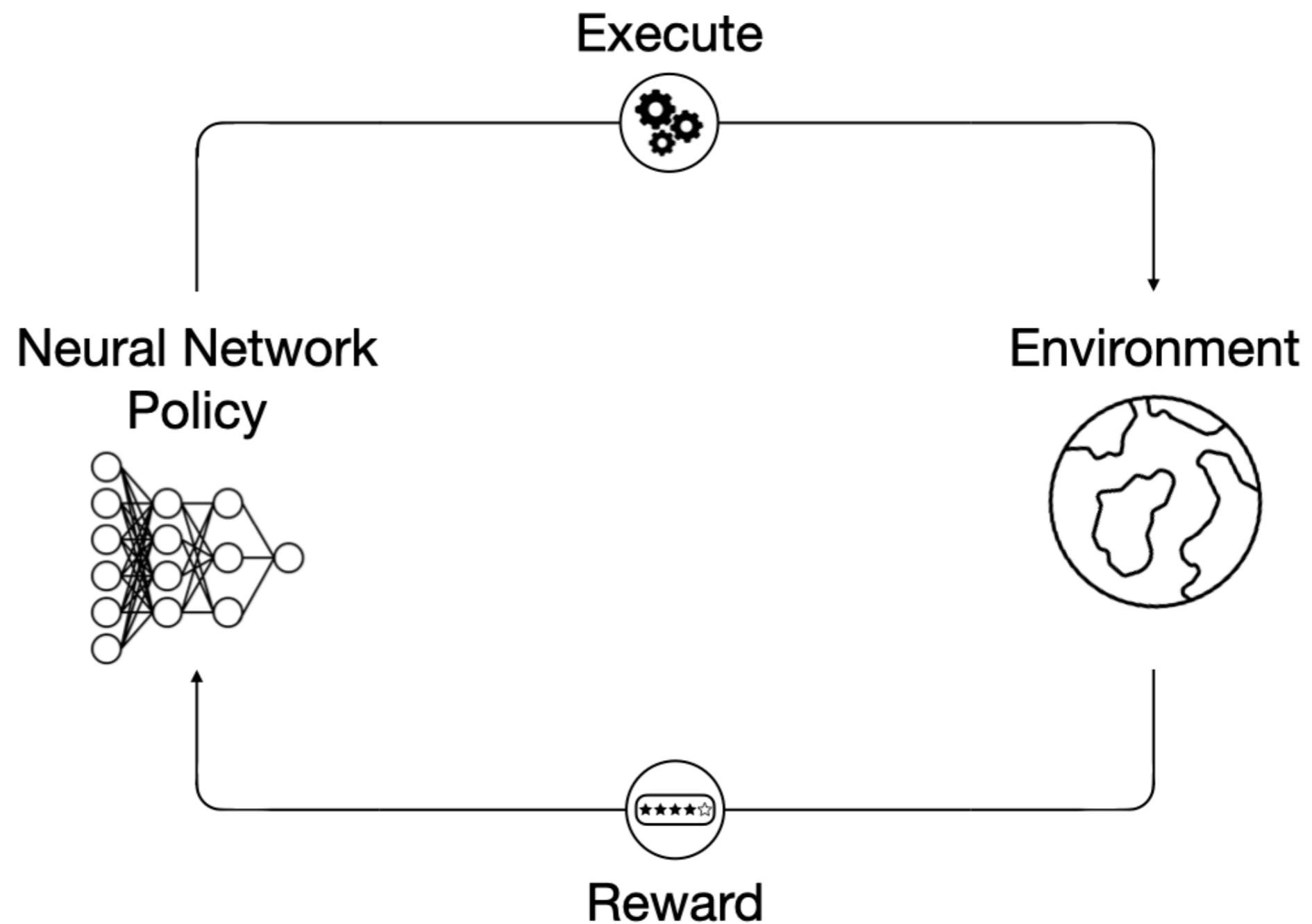


Learning to Synthesize Programs as Interpretable and Generalizable Policies

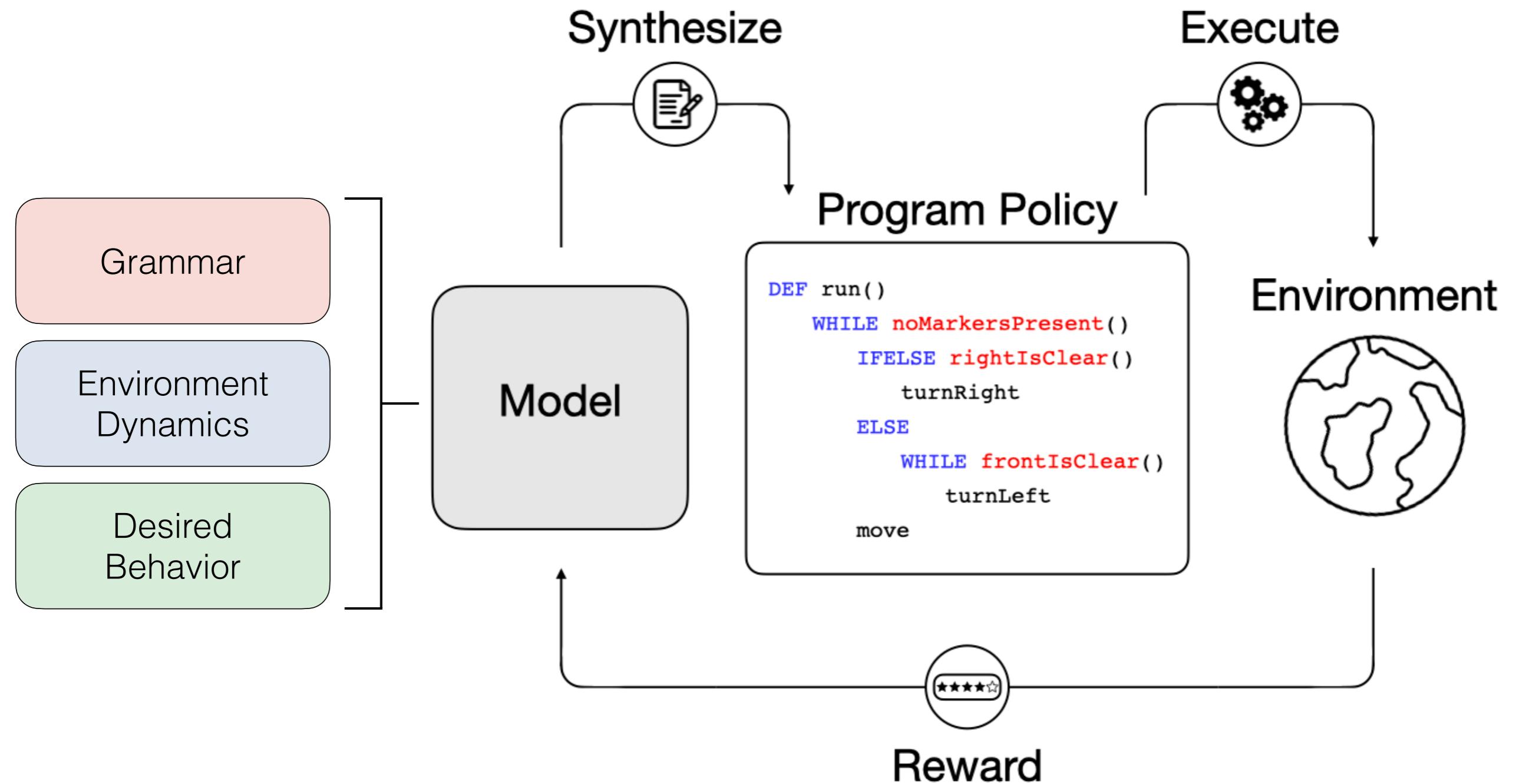
NeurIPS 2021



Reinforcement Learning



Reinforcement Learning by Synthesizing Programs



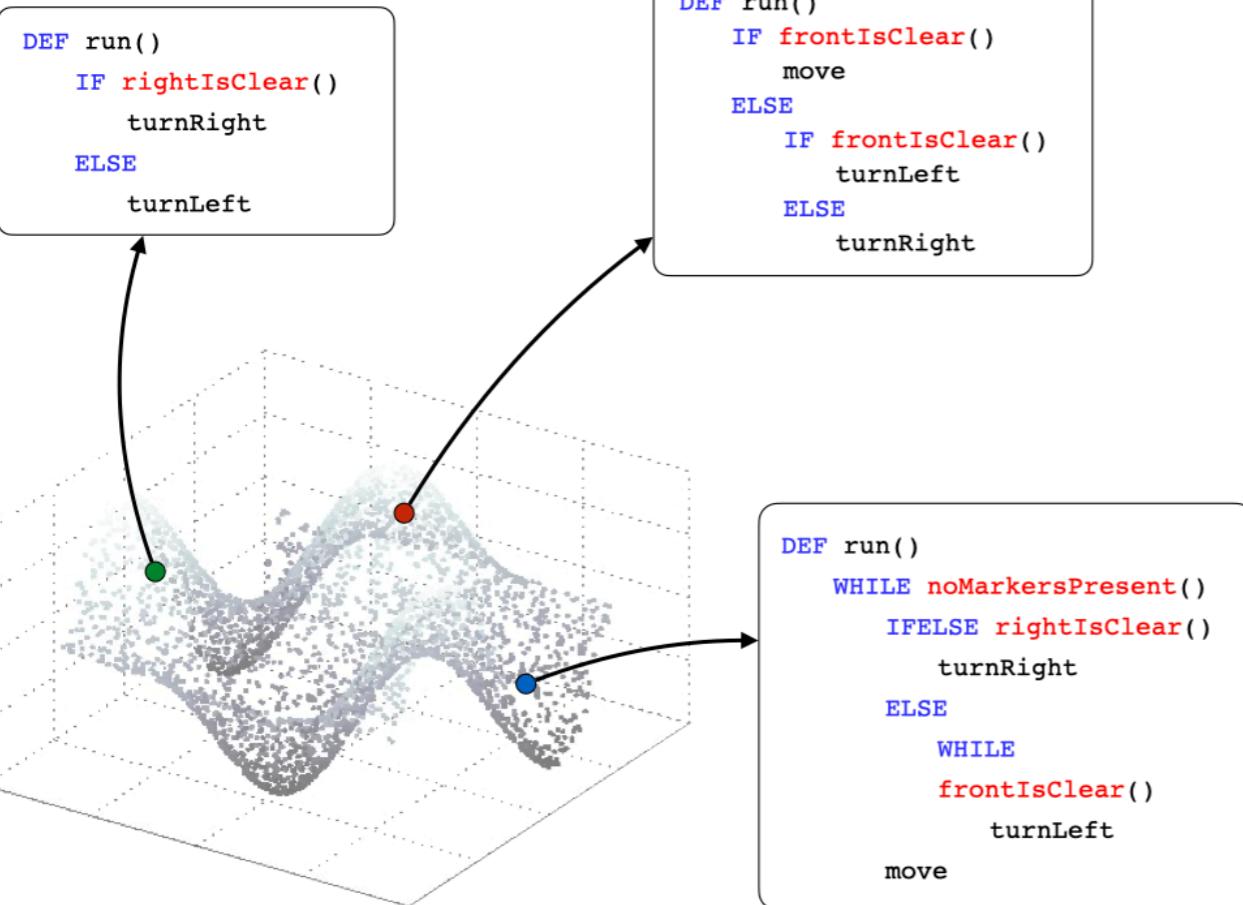
LEAPS: Learning Embeddings for Latent Program Synthesis

Stage 1

Learn a program embedding space from randomly generated programs

Grammar

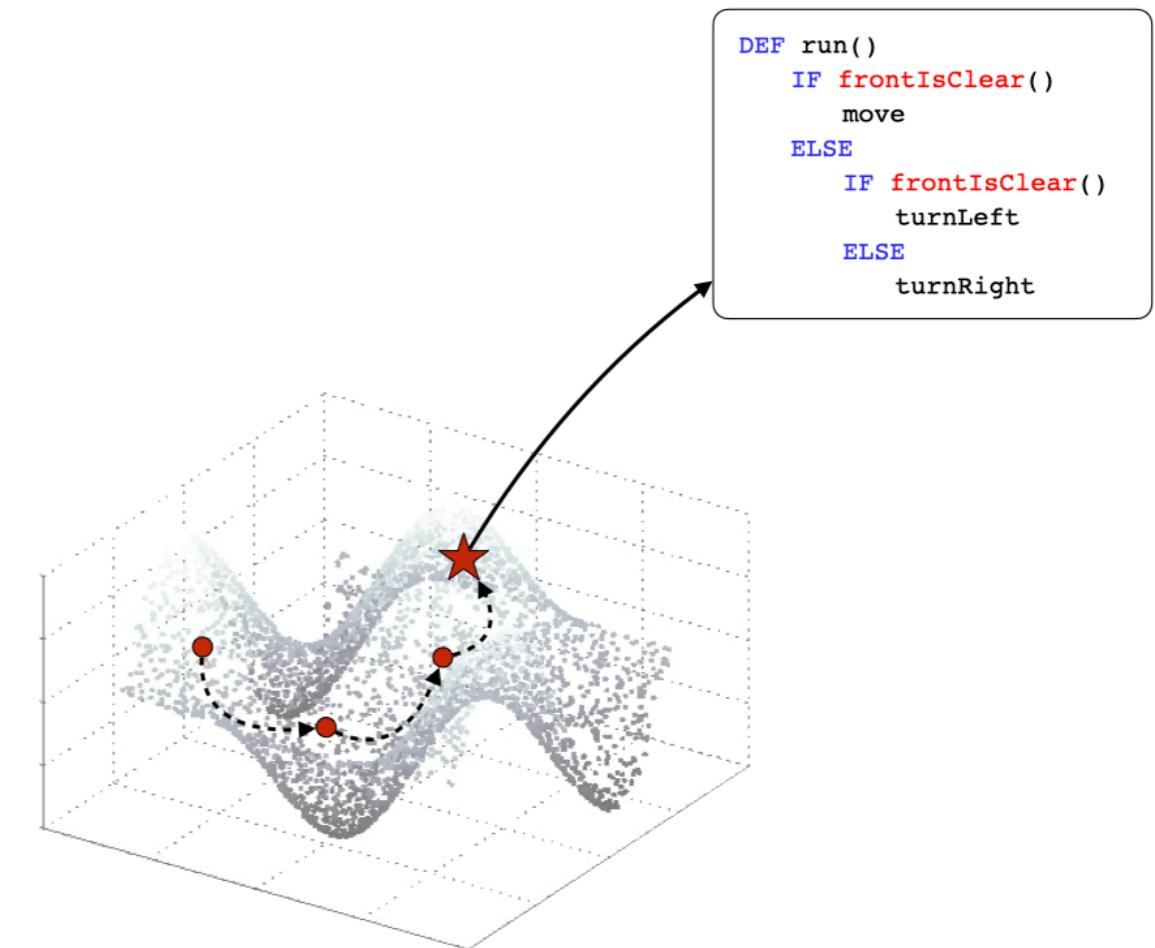
Environment
Dynamics



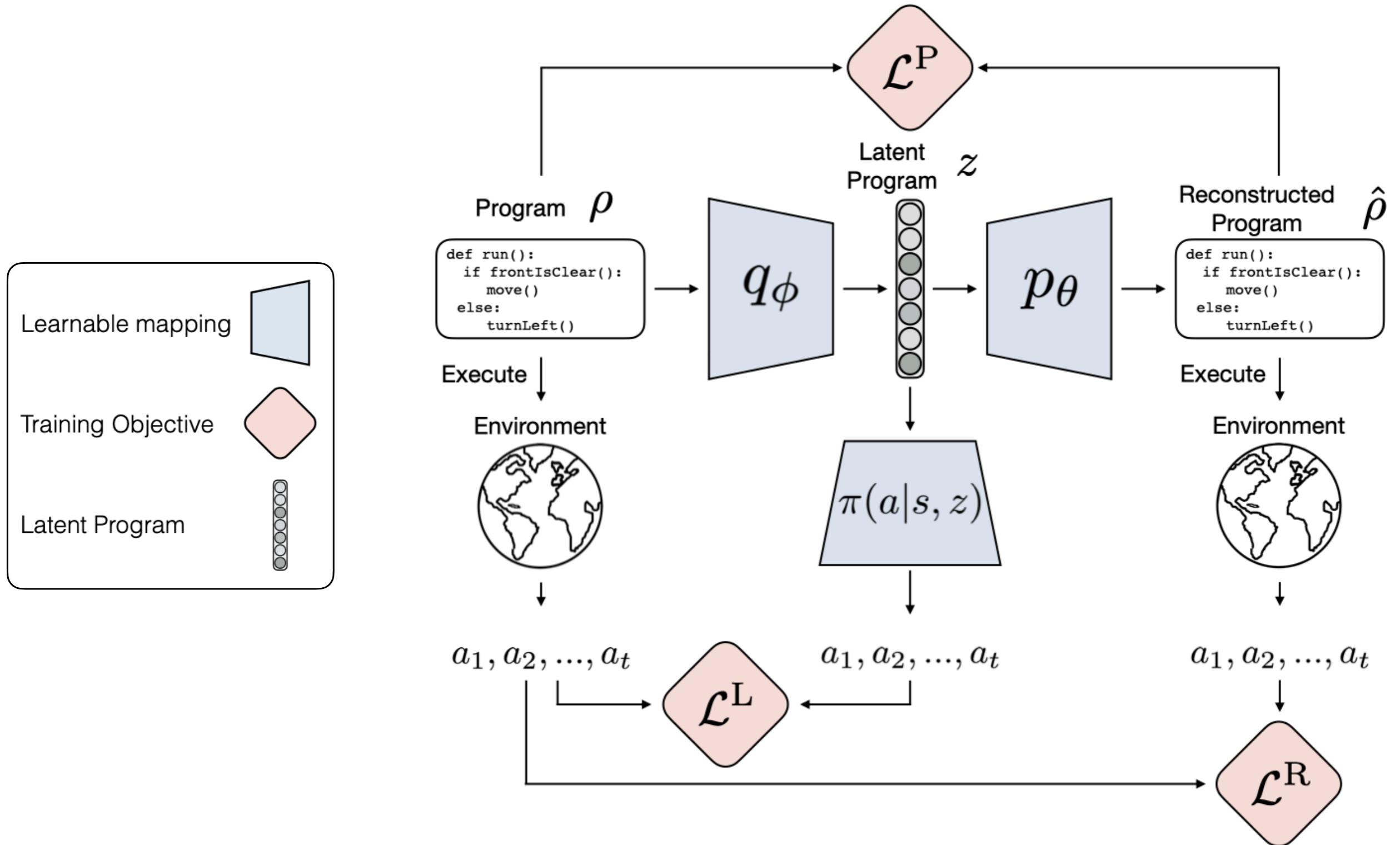
Stage 2

Search for a task-solving program

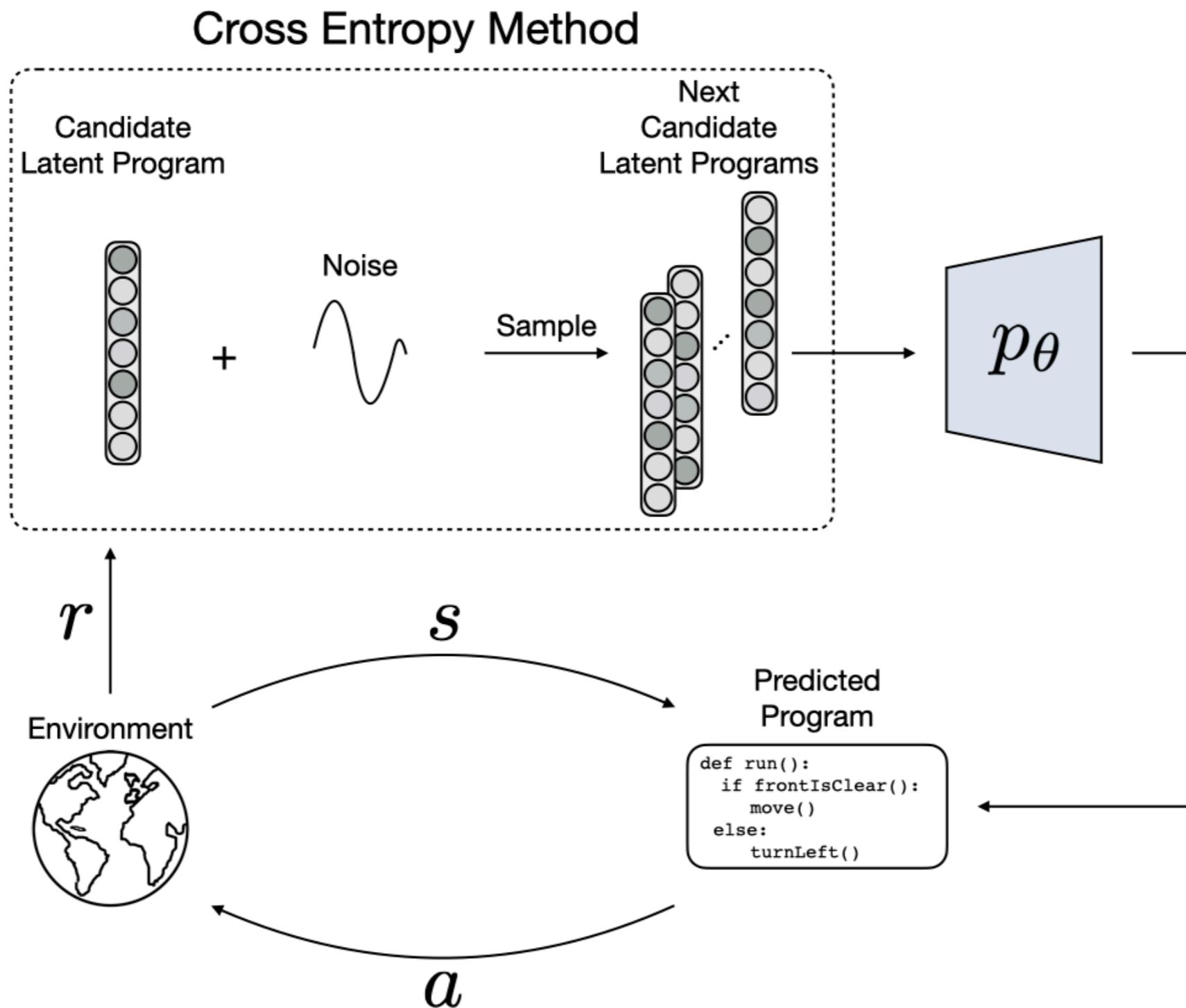
Desired Behavior



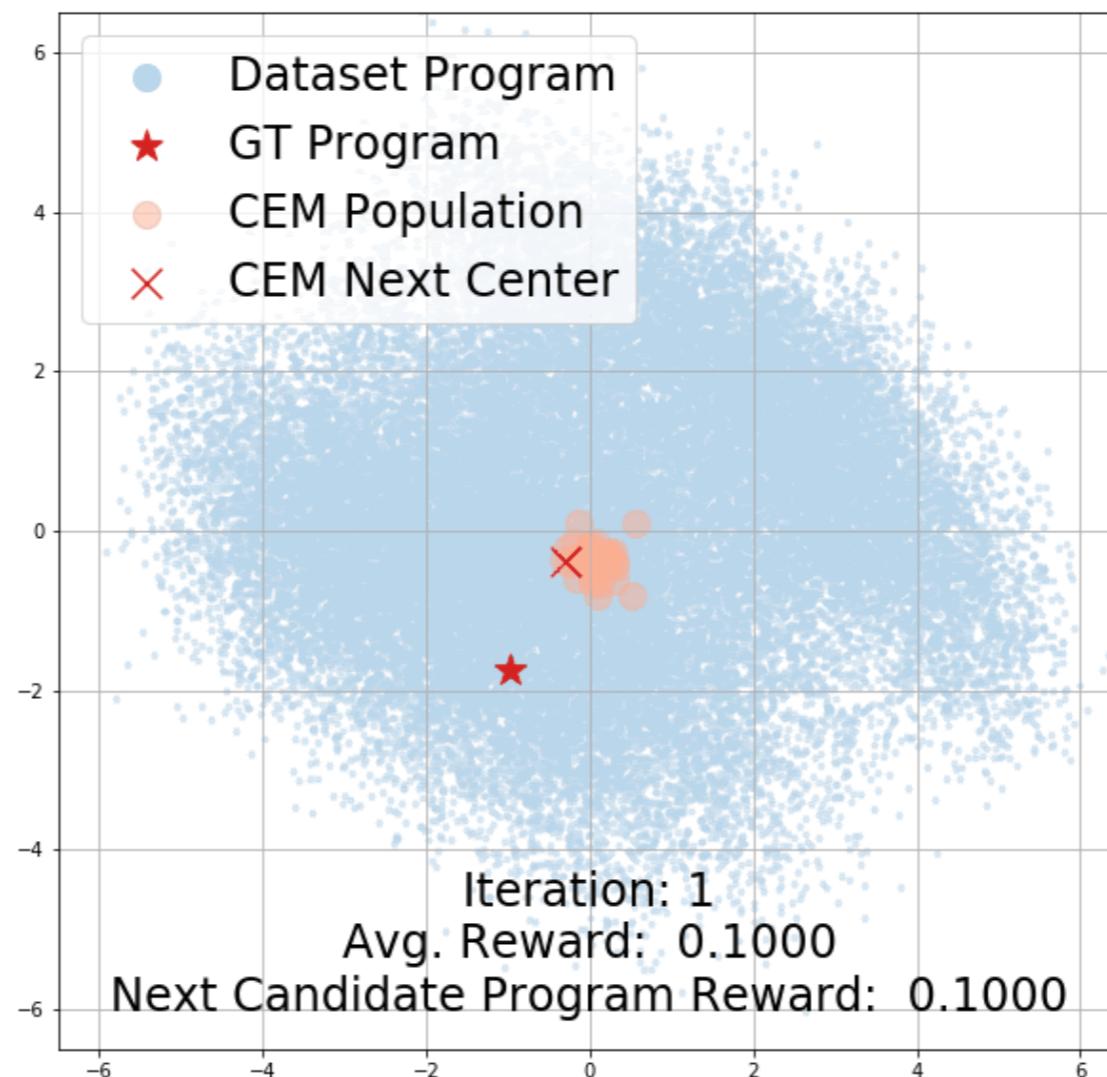
Stage 1: Learning a Program Embedding Space



Stage 2: Latent Program Search with Cross Entropy Method

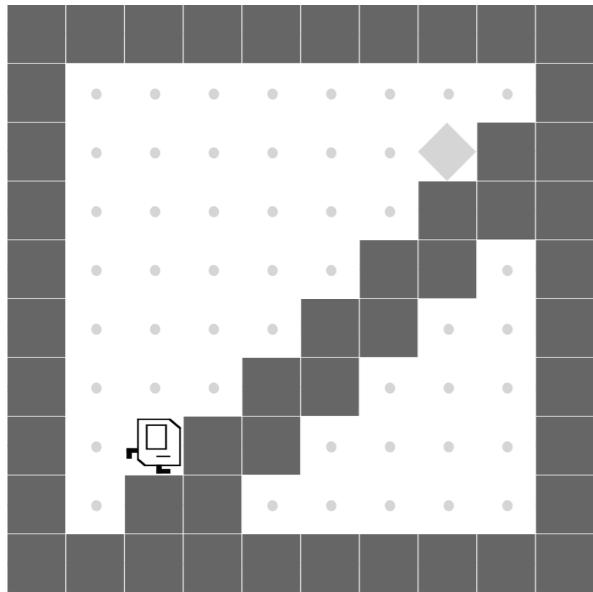


Cross Entropy Method Trajectory Visualization

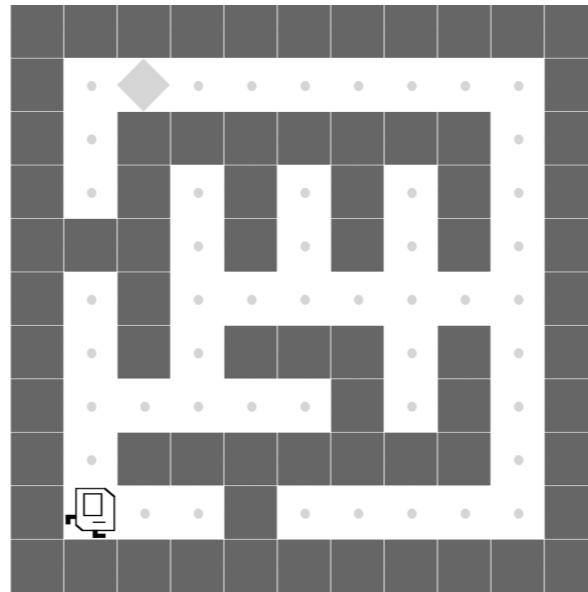


Karel Tasks

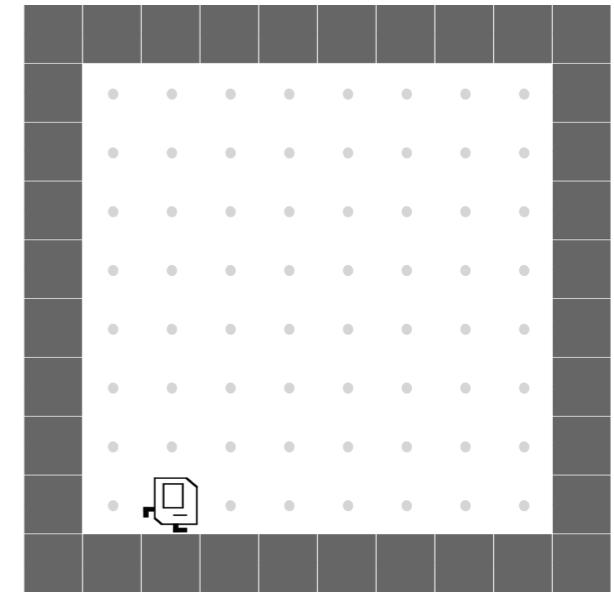
StairClimber



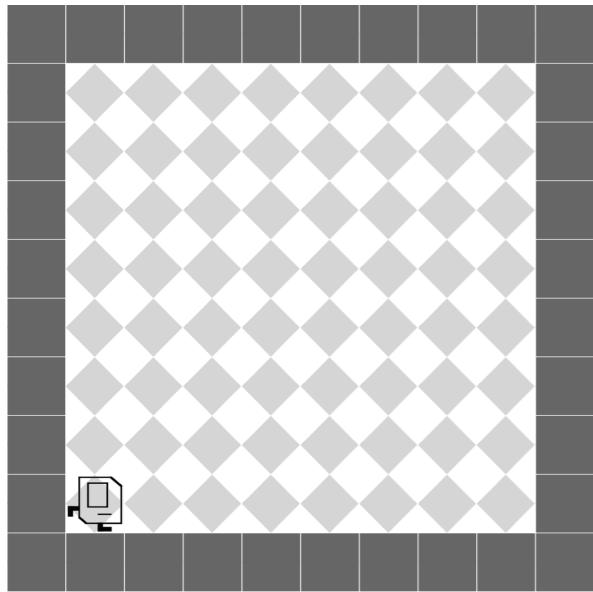
Maze



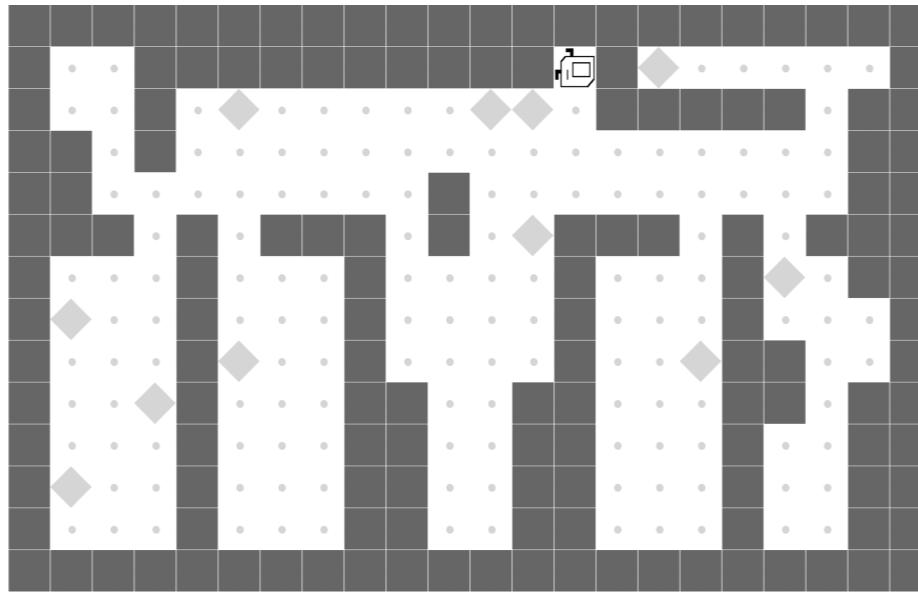
FourCorners



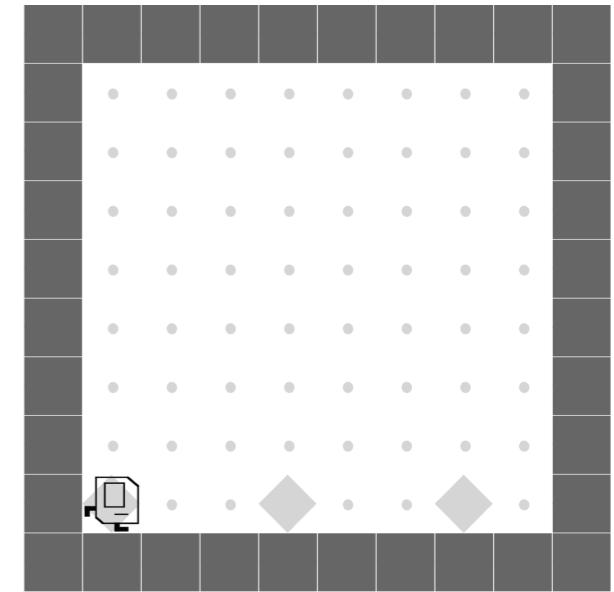
Harvester



CleanHouse

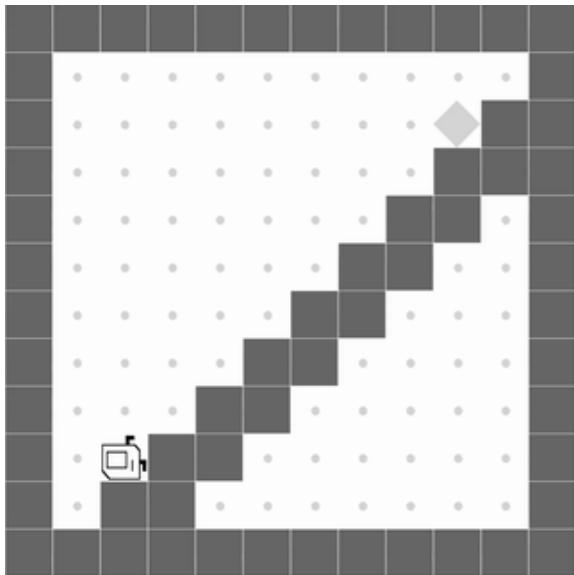


TopOff

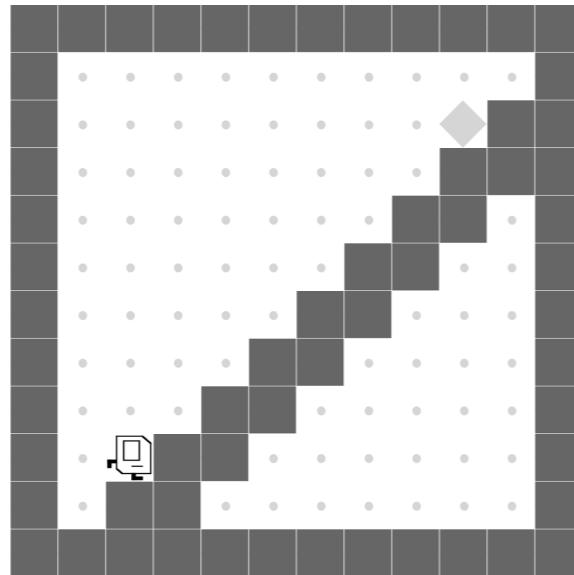


Qualitative Results

StairClimber

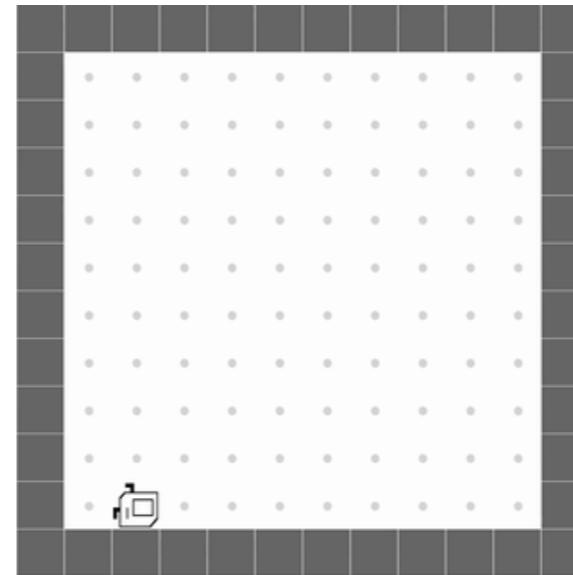


DRL

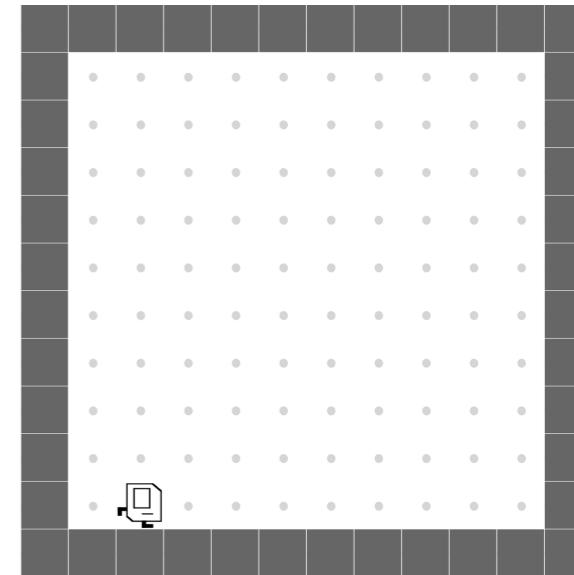


LEAPS

FourCorners

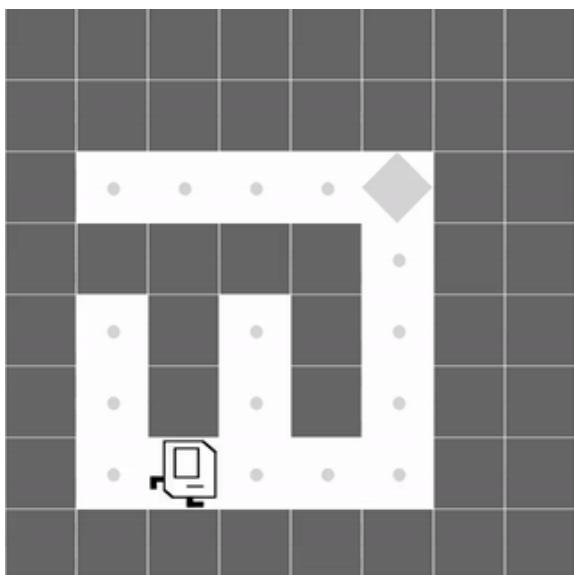


DRL

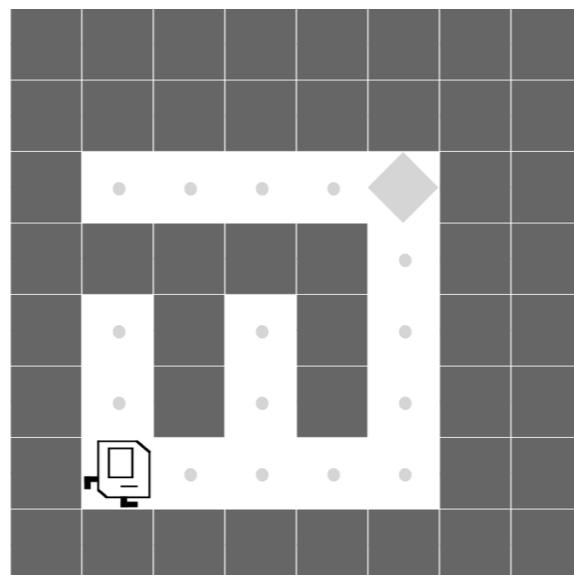


LEAPS

Maze

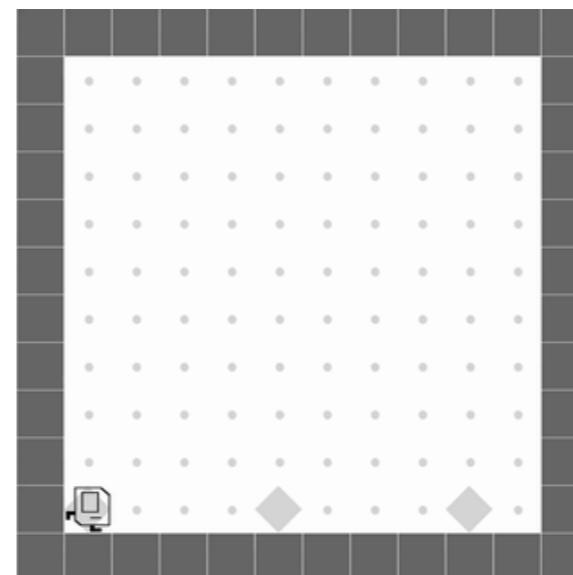


DRL

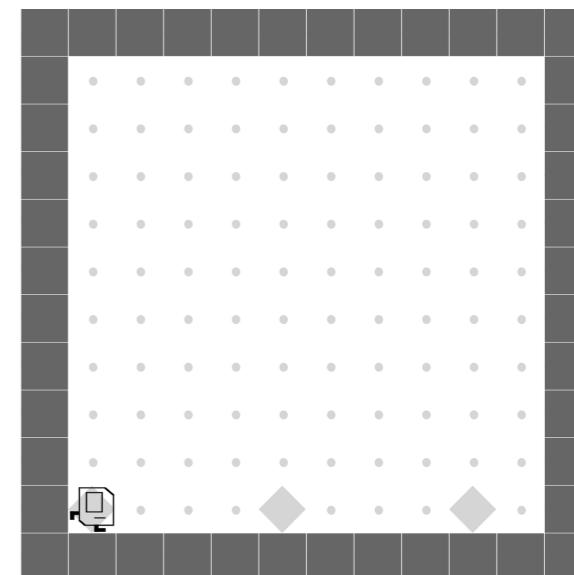


LEAPS

TopOff

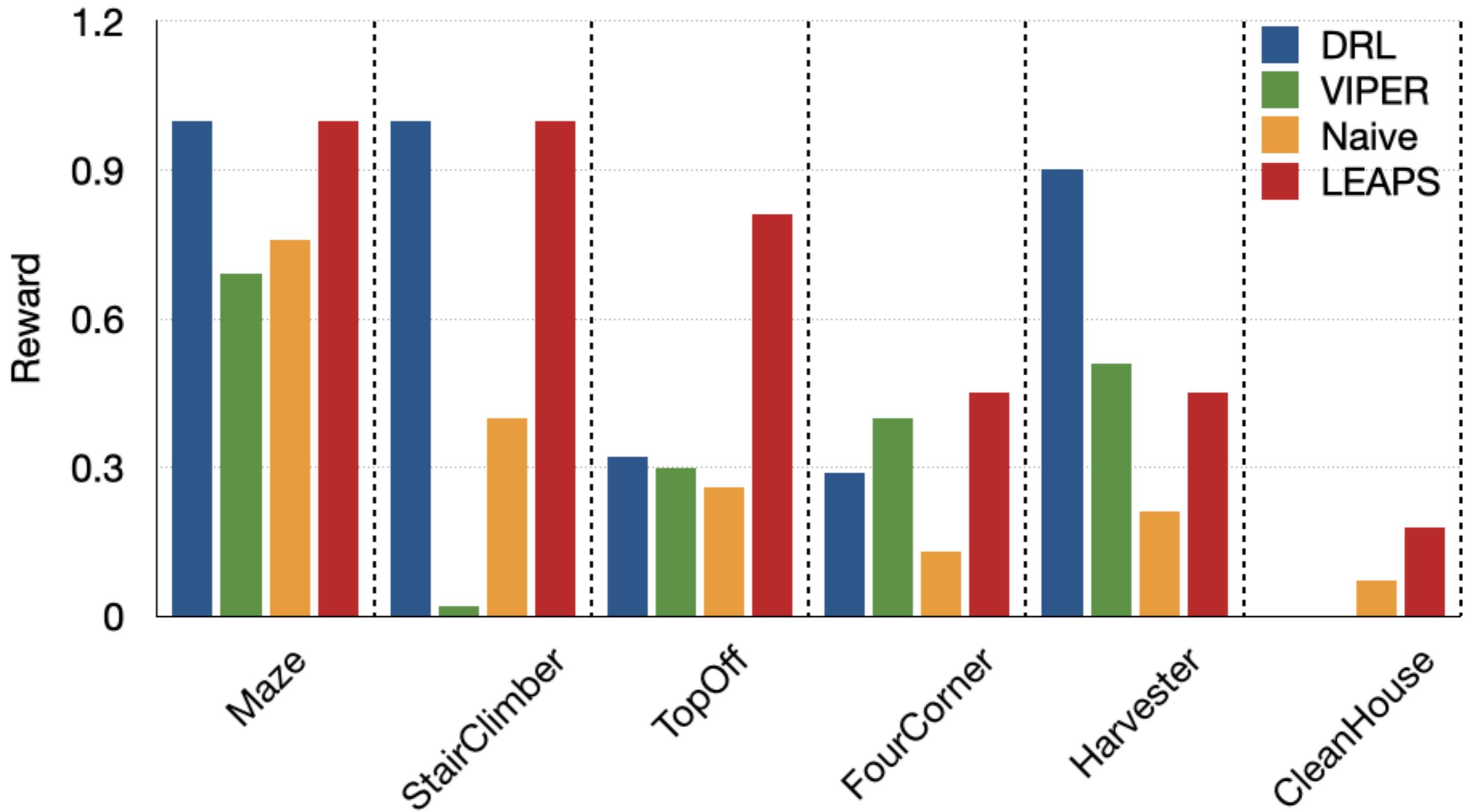


DRL



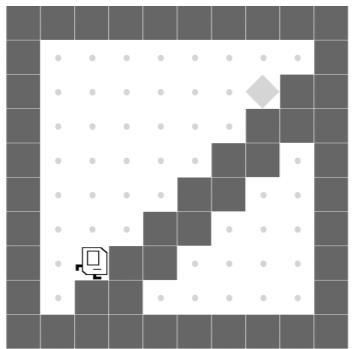
LEAPS

Quantitative Results



Zero-shot Generalization

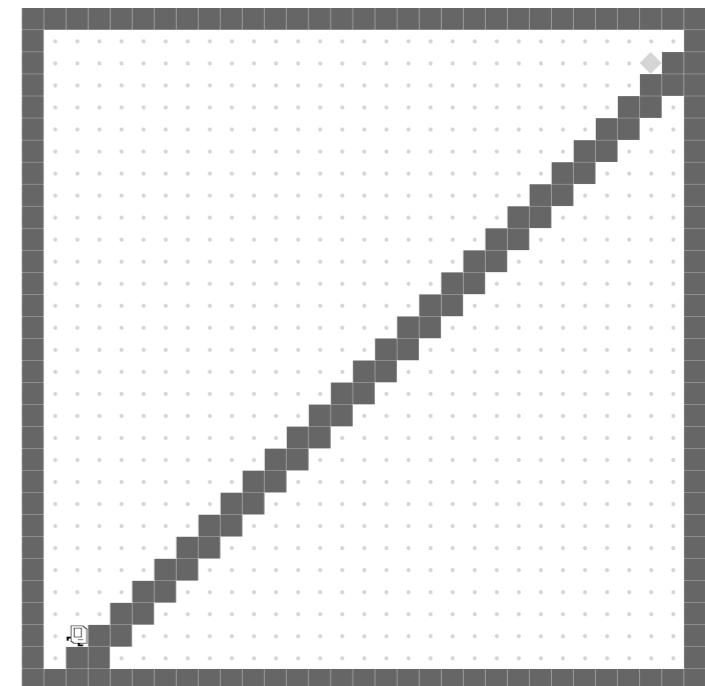
Learning on
8x8 grids



```
DEF run()
  WHILE noMarkersPresent()
    turnRight
    move
  WHILE rightIsClear()
    turnLeft
```

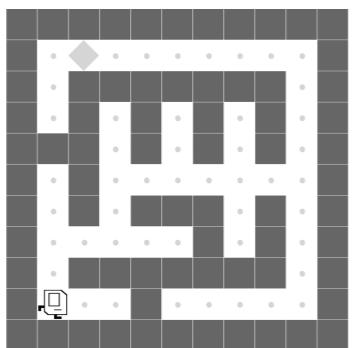


Evaluation on
100x100 grids

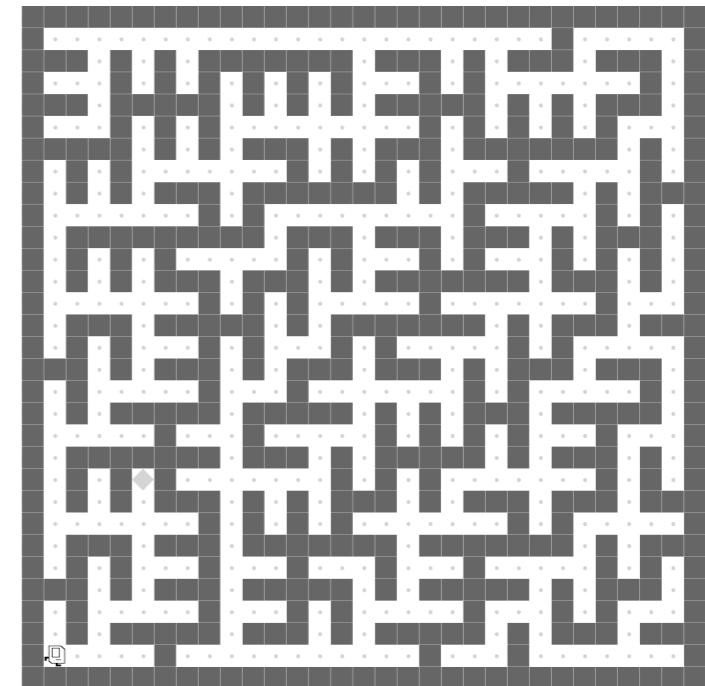


StairClimber

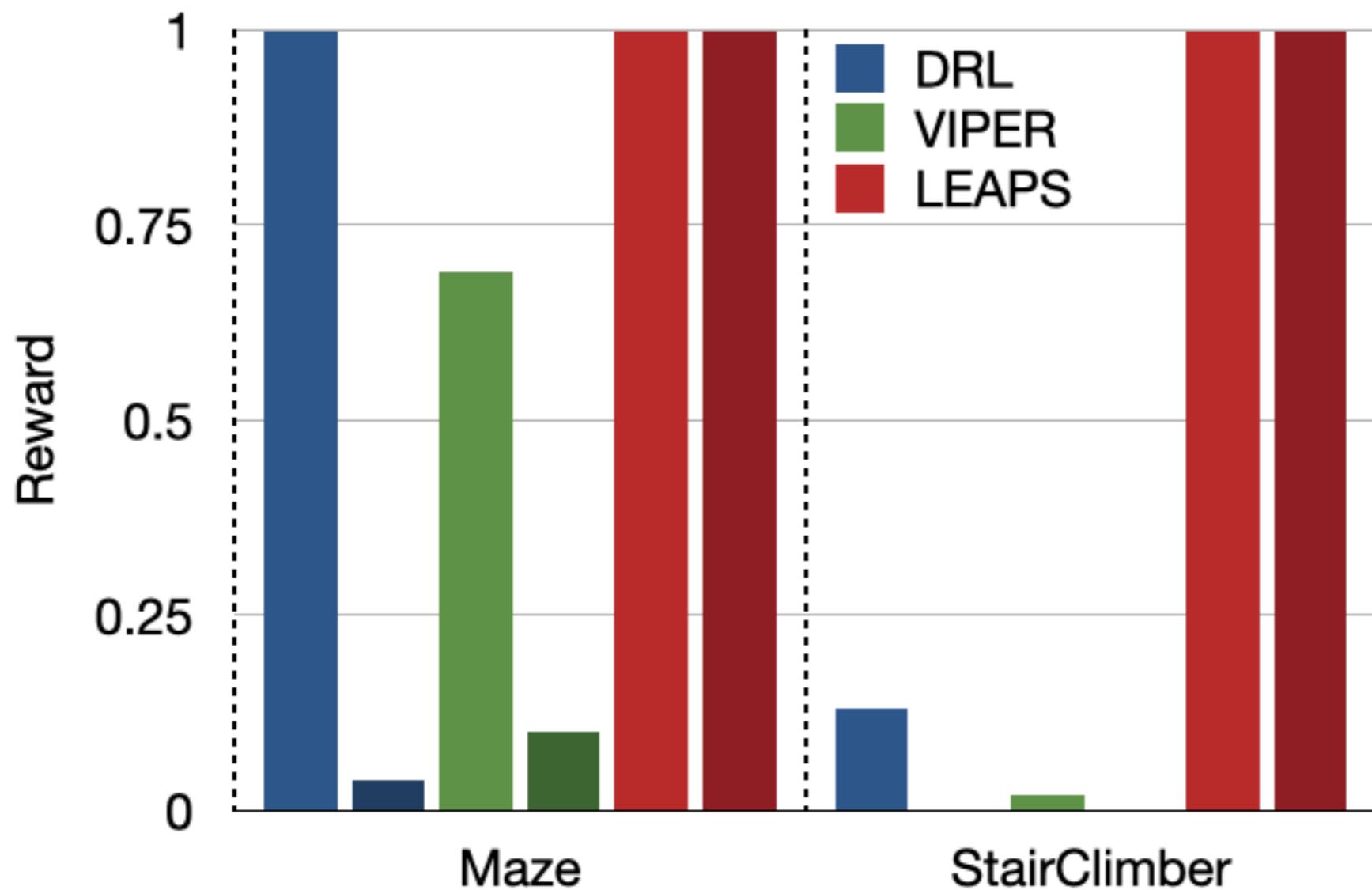
Maze



```
DEF run()
  IF frontIsClear()
    turnLeft
  WHILE noMarkersPresent()
    turnRight
    move
```

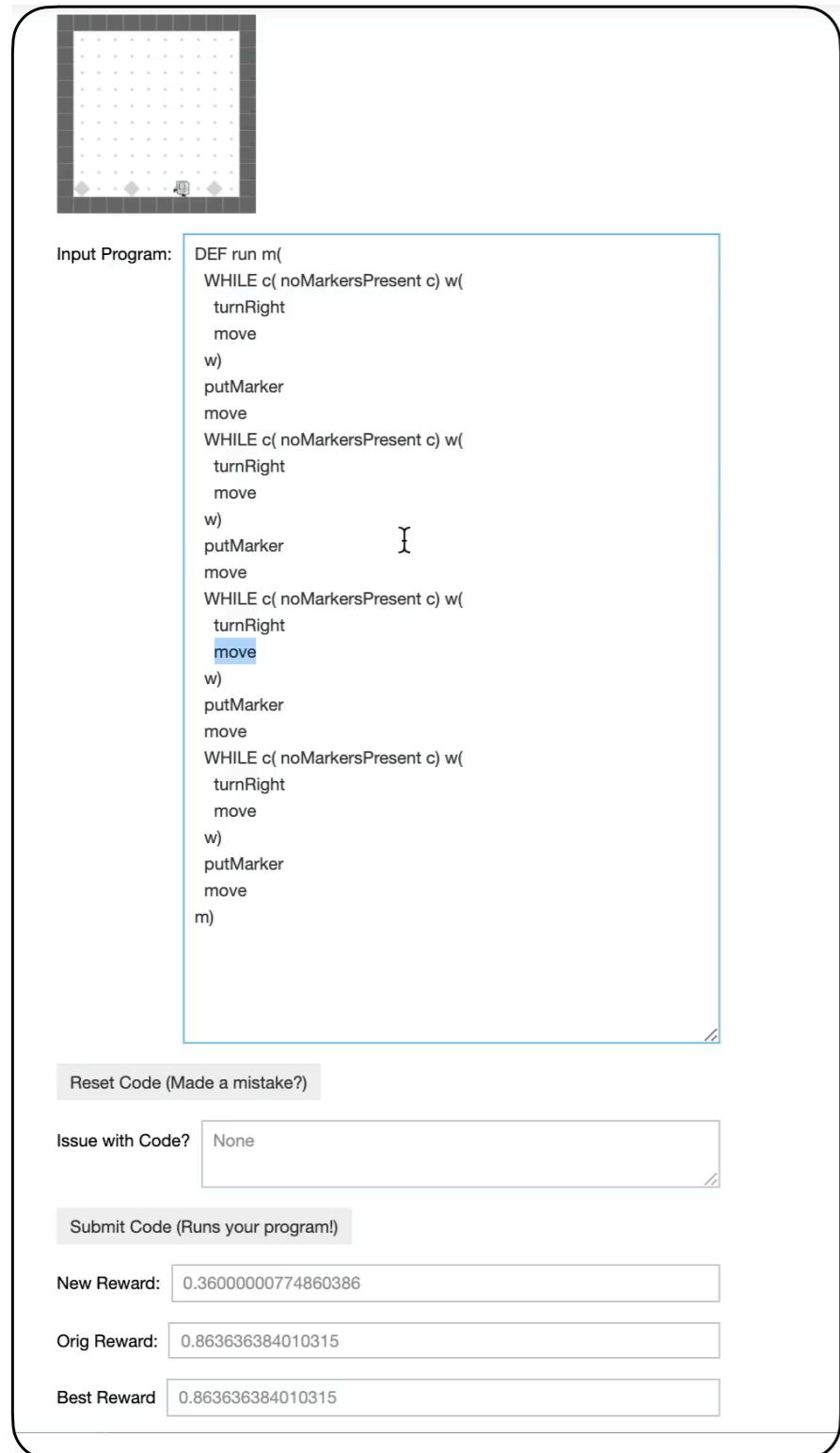


Zero-shot Generalization

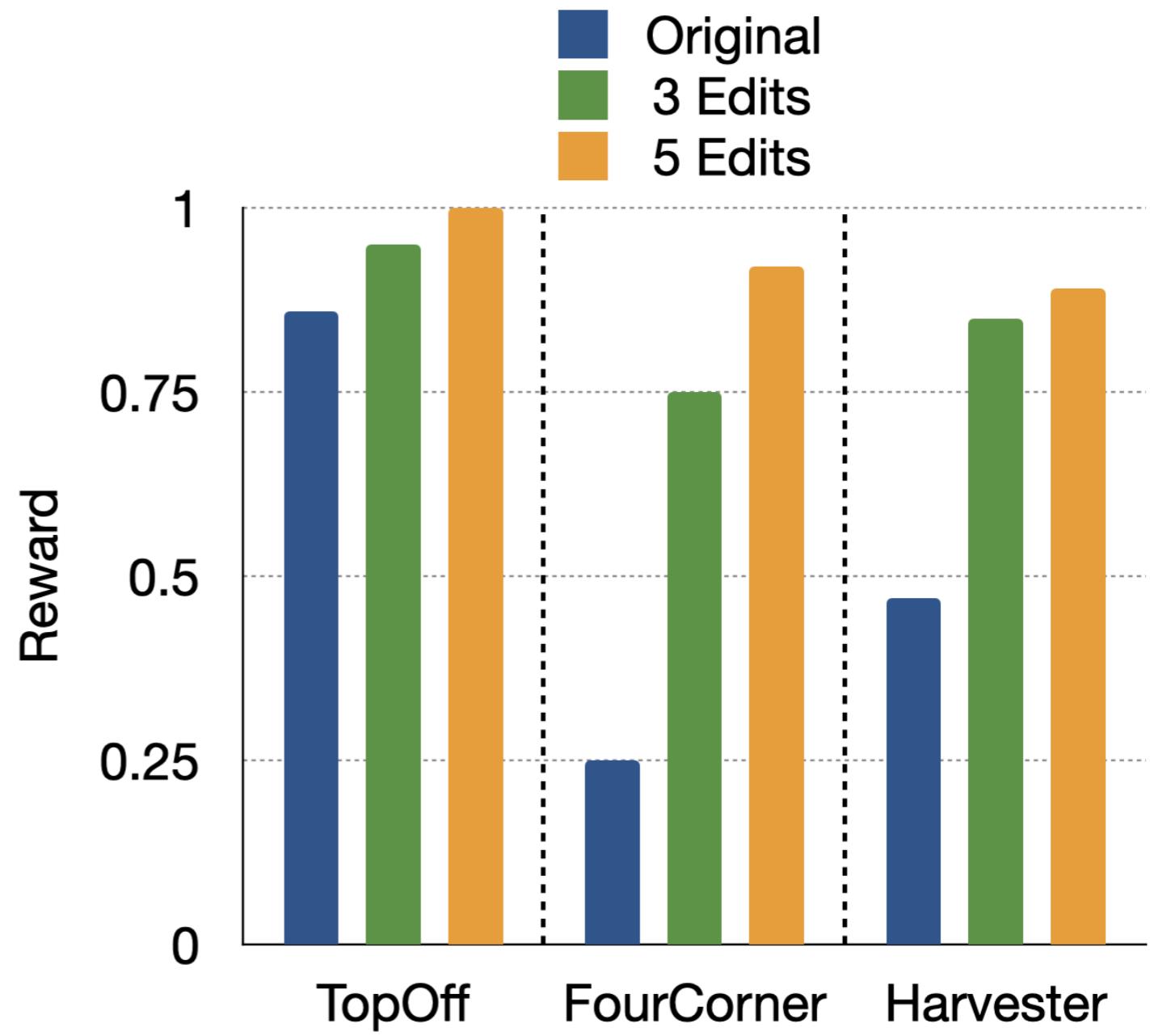


Interpretability

Human Debugging Interface

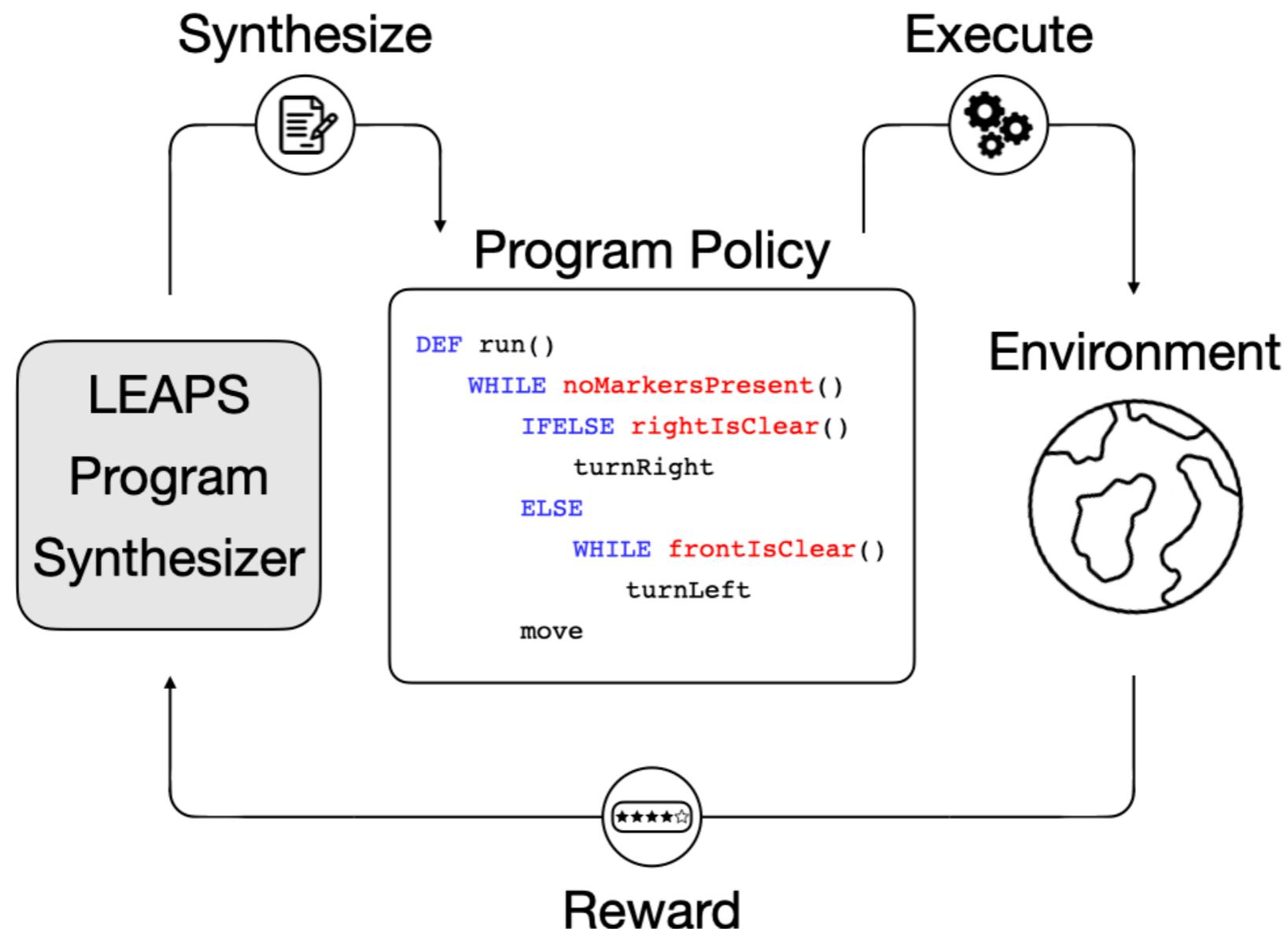


Improved Performance

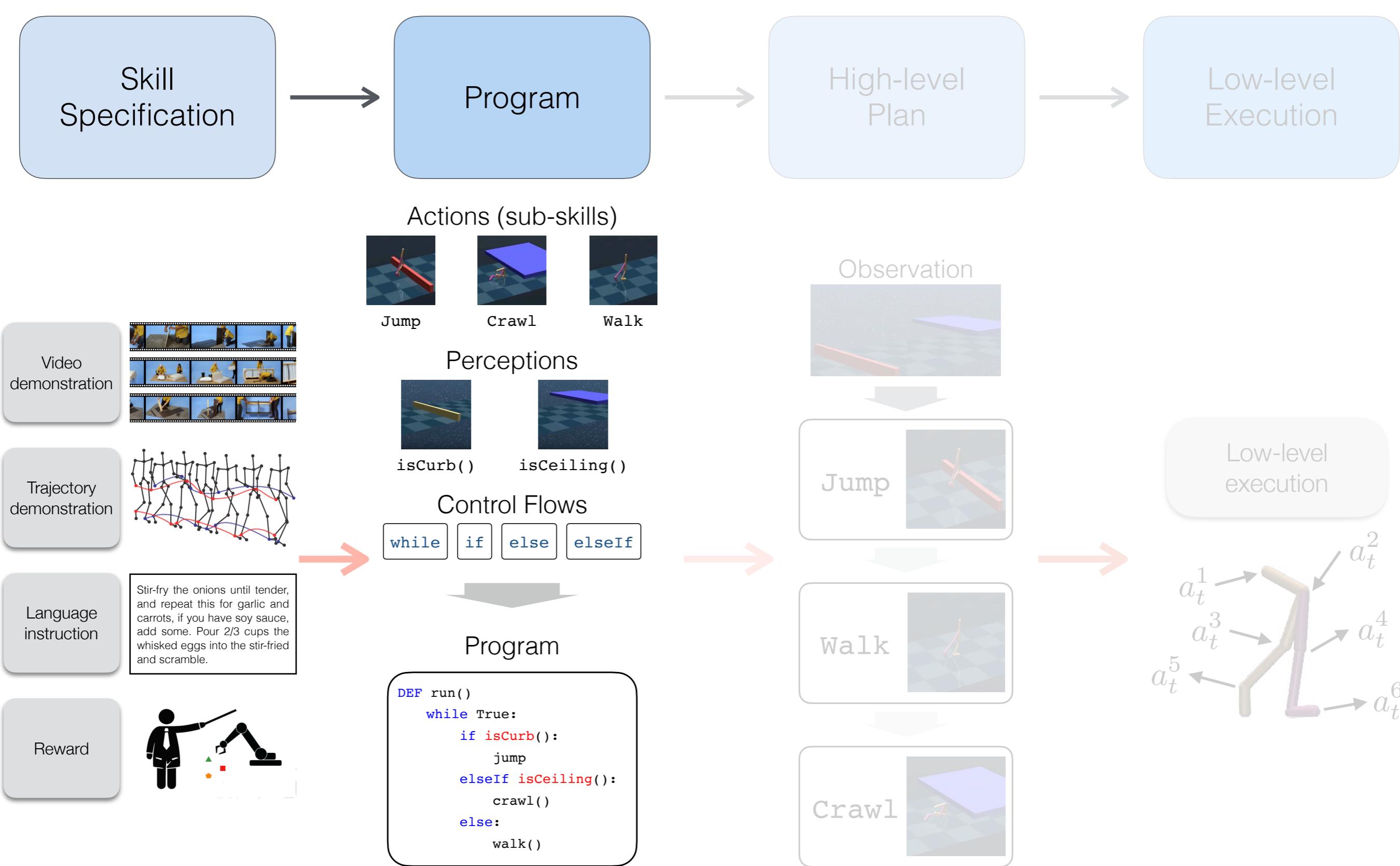


Takeaway

- Synthesize generalizable and interpretable programs from rewards

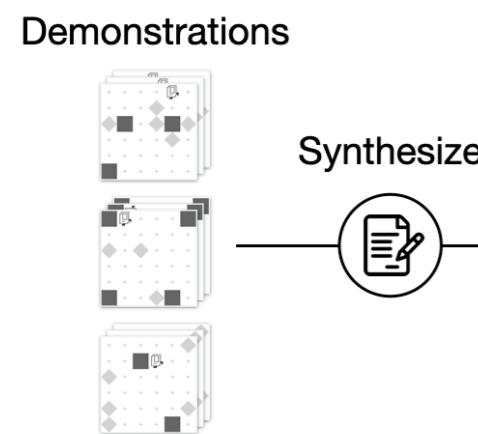


Program Inference

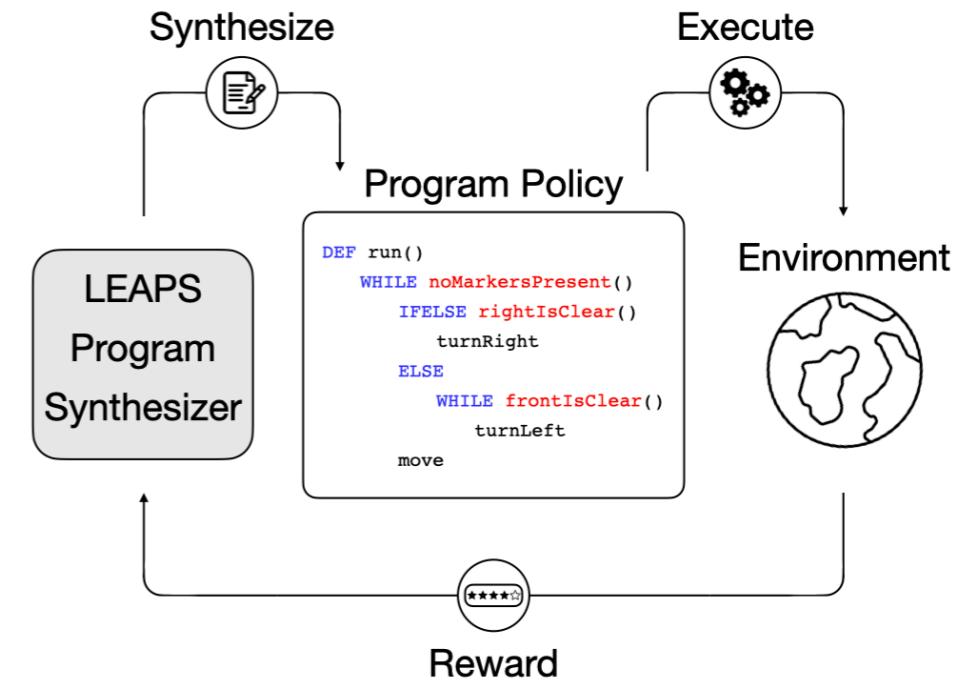


Program Inference

Imitation learning from demonstrations



Reinforcement learning from rewards



Program Inference

Imitation learning from demonstrations

Demonstrations



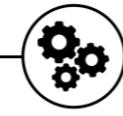
Synthesize



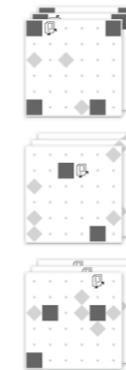
Program

```
DEF run()
  if isFrontClear():
    move
  else:
    turnLeft
    move
    turnLeft
  repeat(2):
    turnRight
    putMarker
```

Execute



Execution



Reinforcement learning from rewards

Synthesize

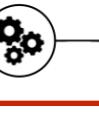


Program Policy

LEAPS
Program
Synthesizer

```
DEF run()
  WHILE noMarkersPresent()
    IFELSE rightIsClear()
      turnRight
    ELSE
      WHILE frontIsClear()
        turnLeft
      move
```

Execute



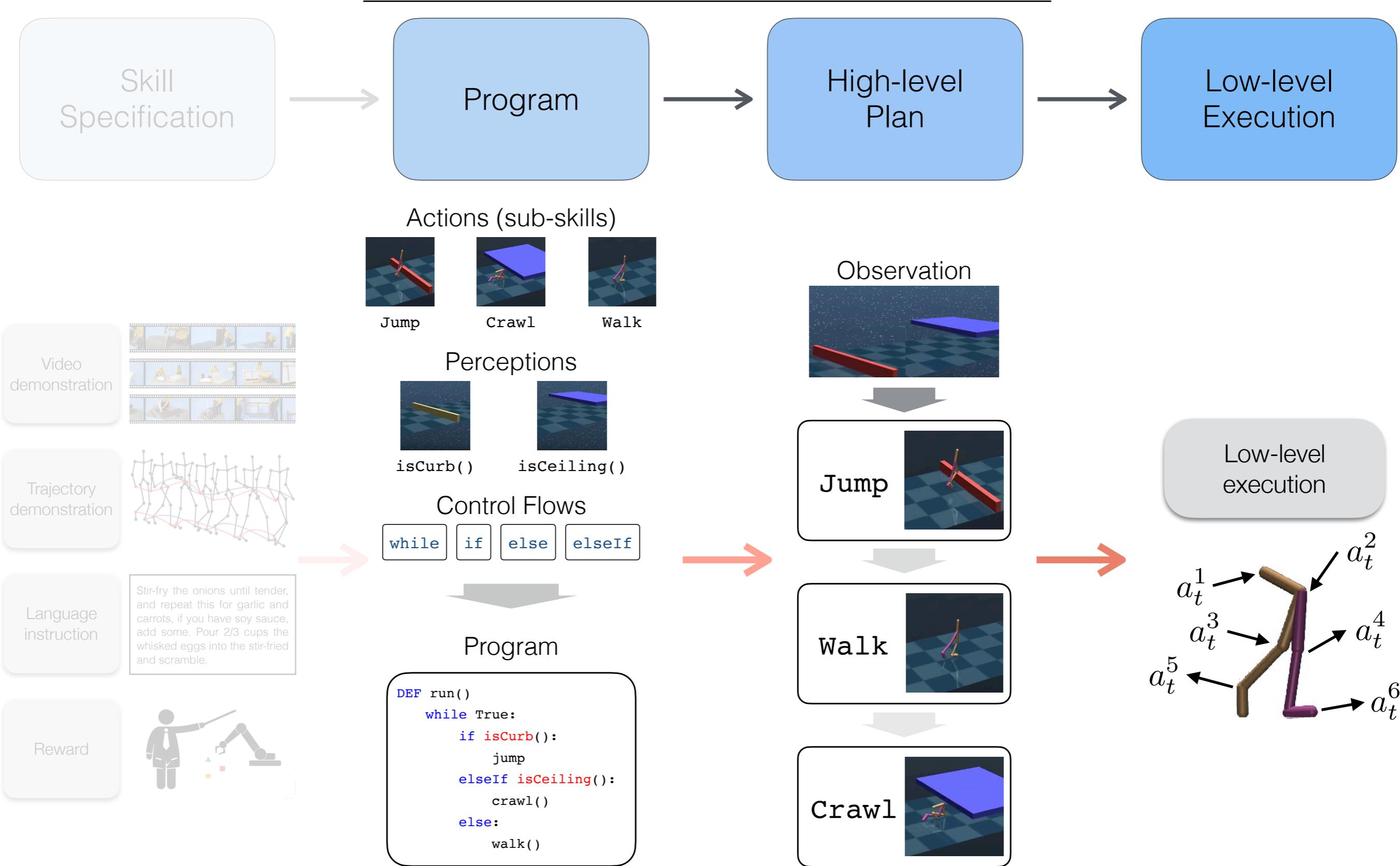
Environment



Reward

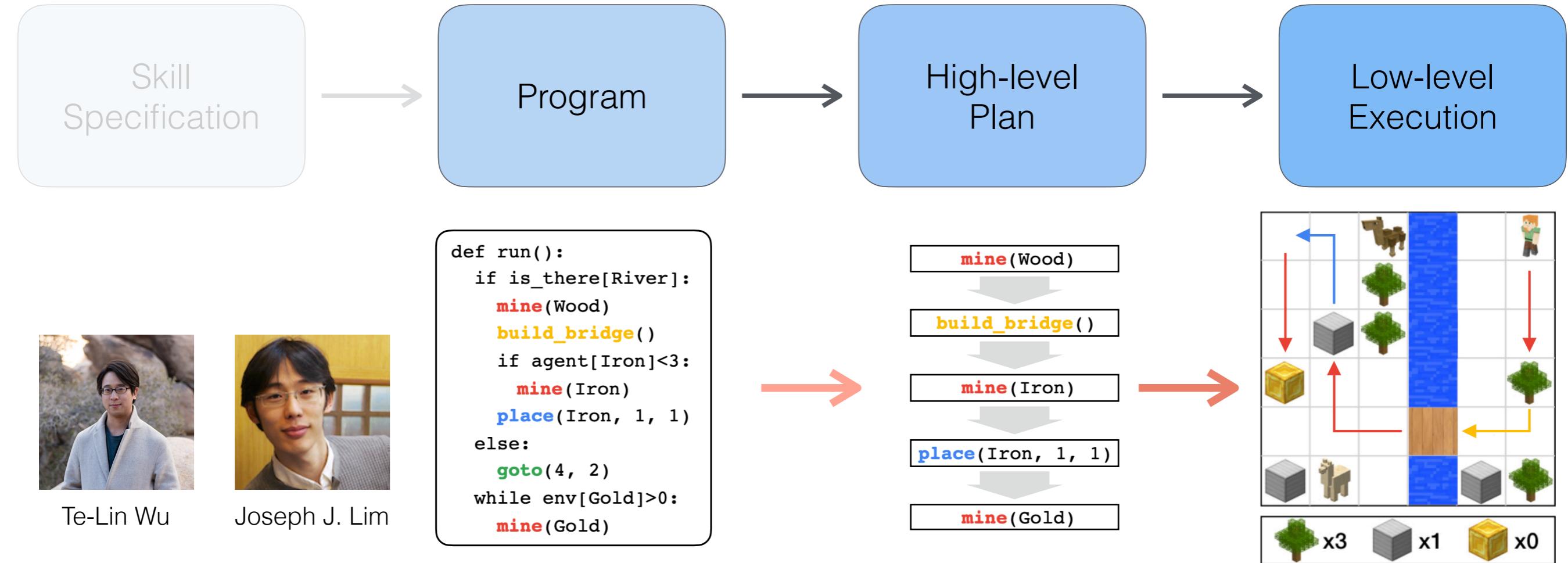


Task Execution



Program Guided Agent

ICLR 2020 (Spotlight)

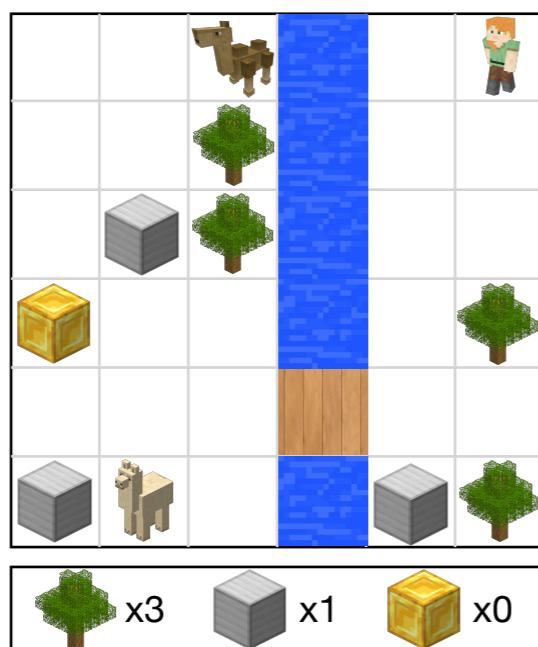


Problem Formulation

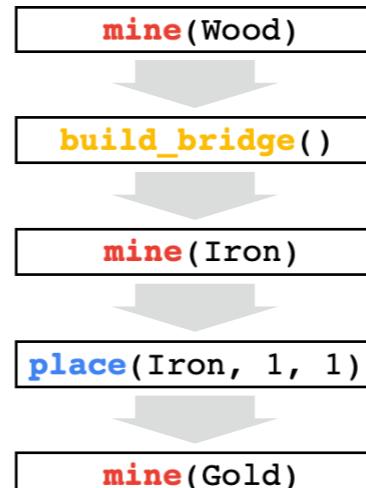
Program
(task)

```
def run():
    if is_there[River]:
        mine(Wood)
        build_bridge()
    if agent[Iron]<3:
        mine(Iron)
        place(Iron, 1, 1)
    else:
        goto(4, 2)
    while env[Gold]>0:
        mine(Gold)
```

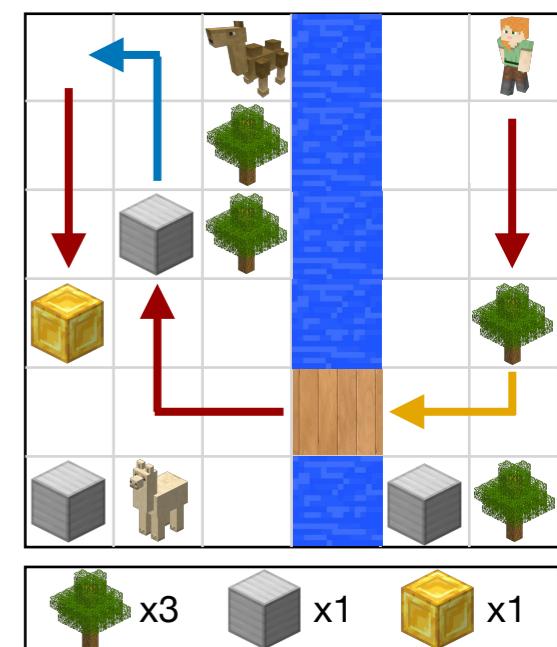
Observation



Plan
(subtasks)



Execution



Instructions

Programs

```
def run()
    if is_there[River]:
        mine(Wood)
        build_bridge()
        if agent[Iron] < 3:
            mine(Iron)
            place(Iron, 2, 3)
        else:
            goto(4, 2)
    while env[Gold] > 0:
        mine(Gold)
```

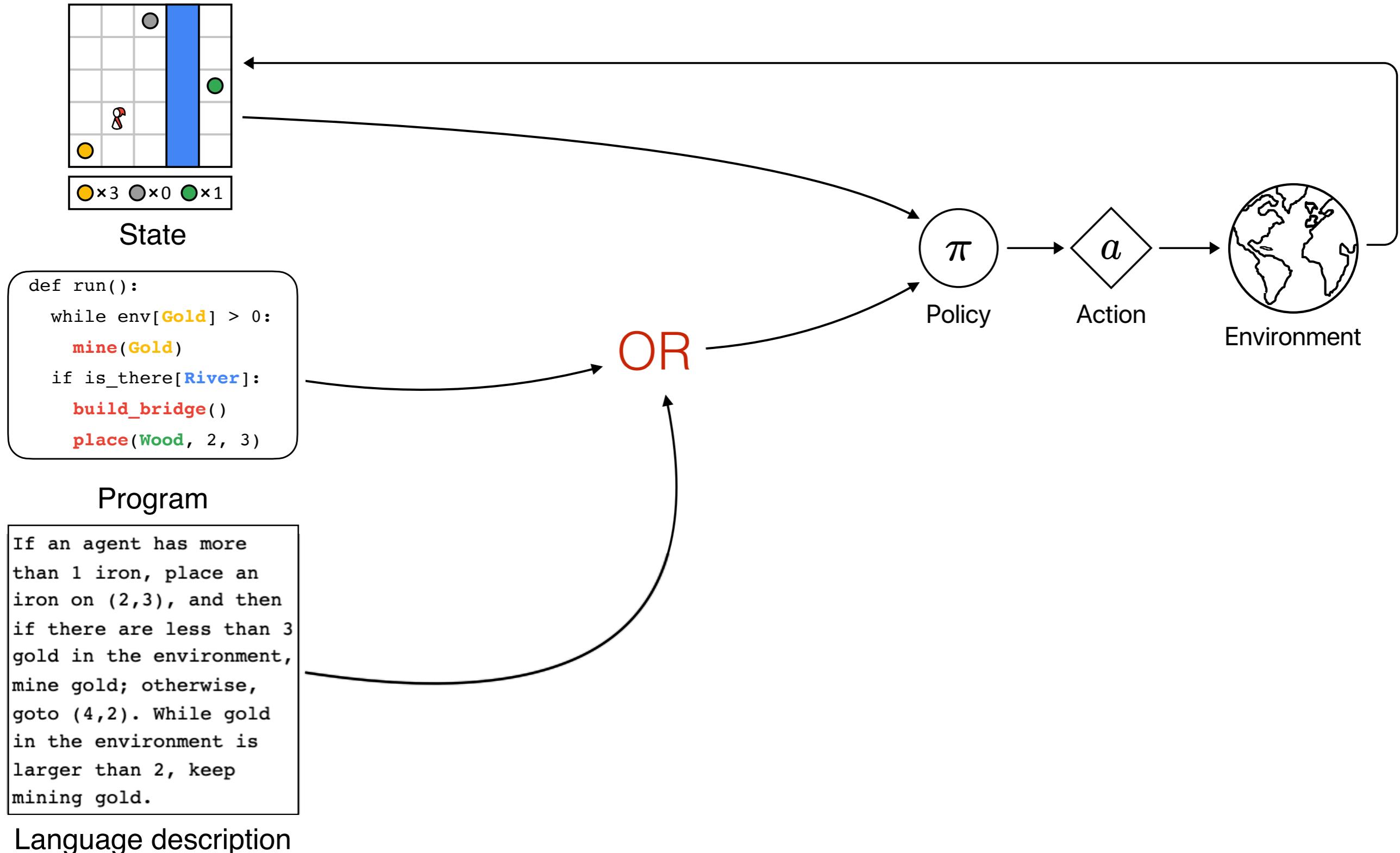
```
def run()
    while agent[Wood] <= 11:
        place(Wood, 2, 4)
        place(Iron, 1, 1)
        place(Iron, 8, 5)
        mine(Gold)
        mine(Gold)
        mine(Gold)
        repeat(4):
            sell(Gold)
            sell(Iron)
```

Natural Language Descriptions

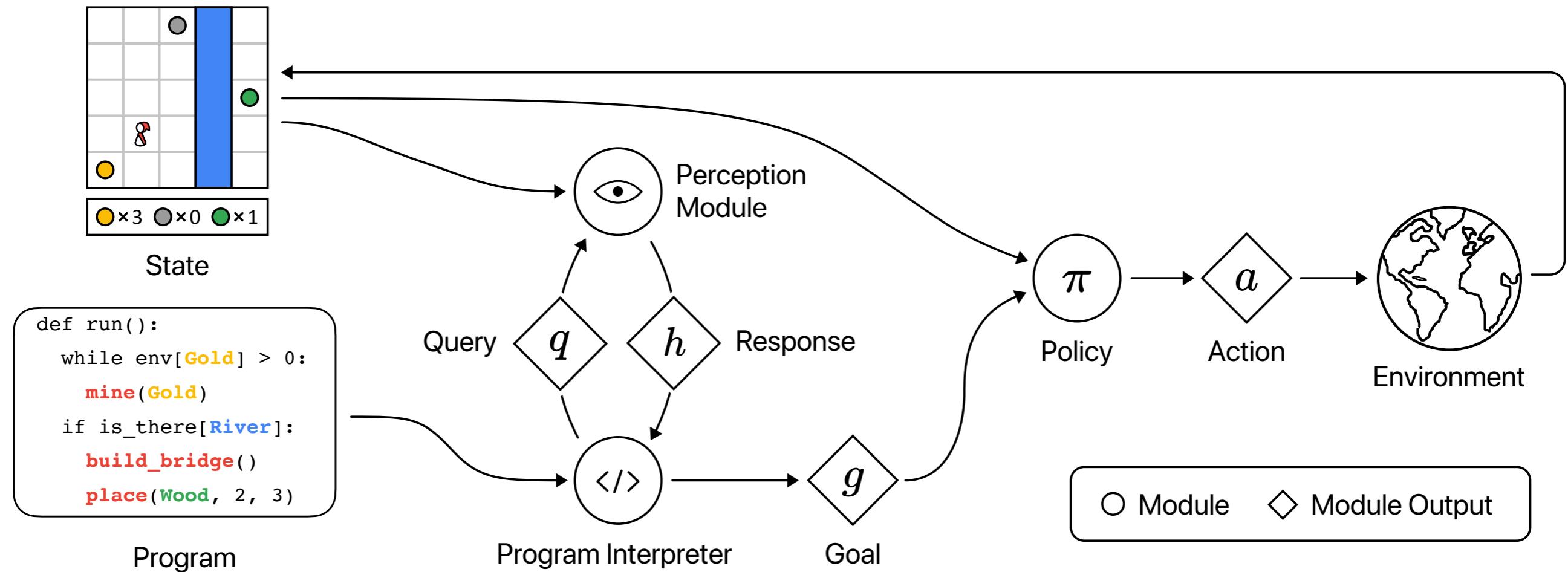
If a river is in the environment, mine a wood and then use it to build a bridge. And then if agent has less than there iron, place an iron at (2,3). Otherwise if no river, goto location (4,2). Finally, whenever there's still gold in the environment, mine a gold.

While agent has no more than 11 wood, place wood at (2,4) and iron at (1,1), then place iron at (8,5) and mine gold twice, then mine gold. After the preceding procedure, sell gold and sell iron 4 times.

End-to-end Learning Baseline



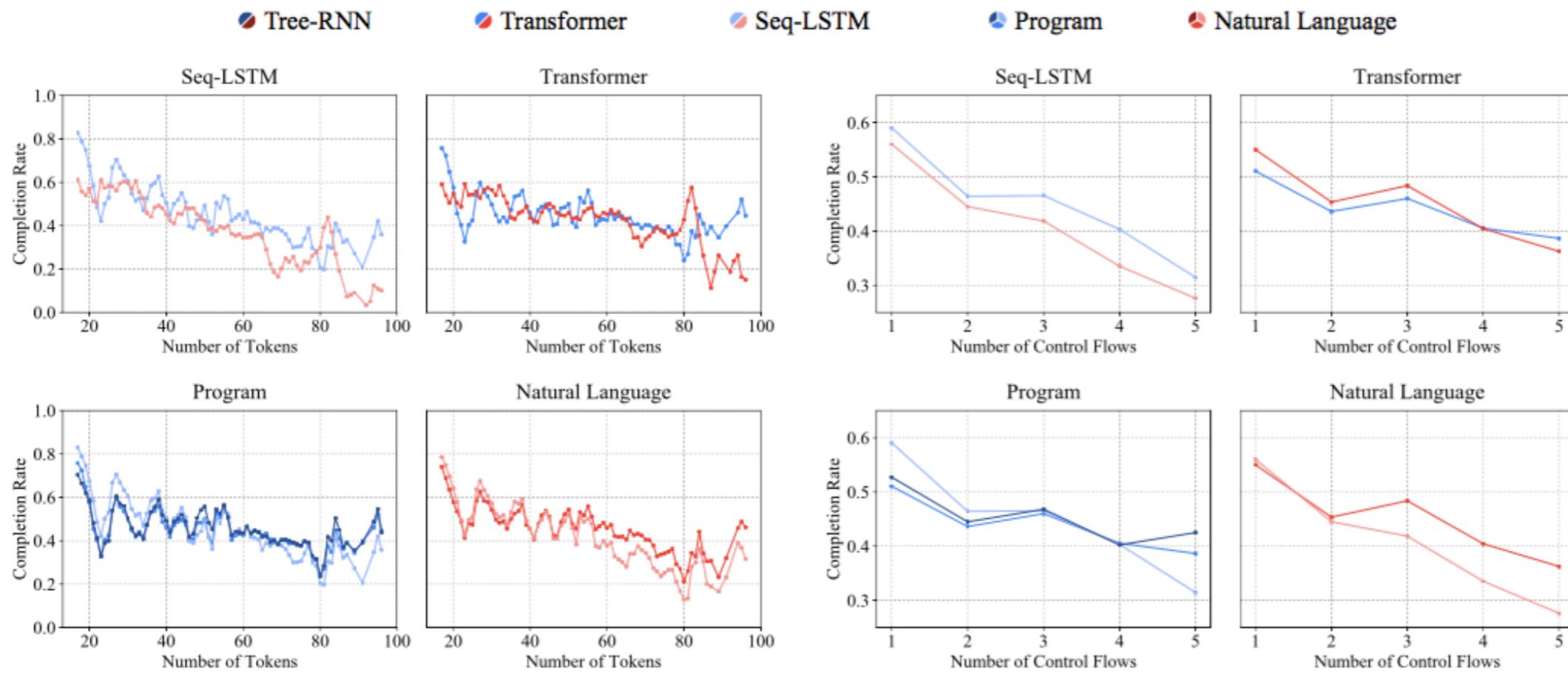
Program Guided Agent



Quantitative Results

Generalization

Instruction Method		Natural language descriptions		Programs				
		Seq-LSTM	Transformer	Seq-LSTM	Tree-RNN	Transformer	Ours (concat)	Ours
Dataset	test	54.9±1.8%	52.5±2.6%	56.7±1.9%	50.1±1.2%	49.4±1.6%	88.6±0.8%	94.0±0.5%
	test-complex	32.4±4.9%	38.2±2.6%	38.8±1.2%	42.2±2.4%	40.9±1.5%	85.2±0.8%	91.8±0.2%
Generalization gap		40.9%	27.2%	31.6%	15.8%	17.2%	3.8%	2.3%



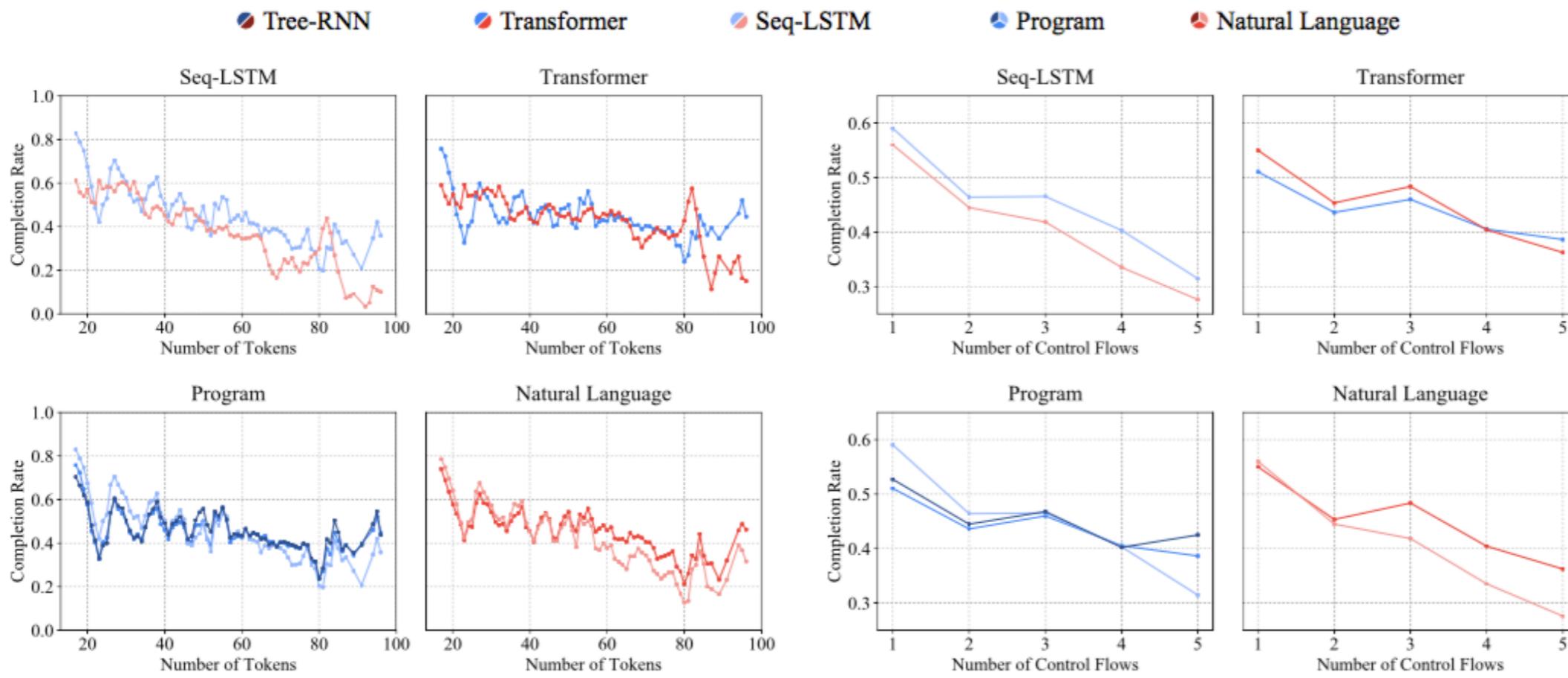
(a) Instruction Length

(b) Instruction Complexity

Quantitative Results

Natural Languages < Programs

Instruction Method		Natural language descriptions		Programs				
		Seq-LSTM	Transformer	Seq-LSTM	Tree-RNN	Transformer	Ours (concat)	Ours
Dataset	test	54.9±1.8%	52.5±2.6%	56.7±1.9%	50.1±1.2%	49.4±1.6%	88.6±0.8%	94.0±0.5%
	test-complex	32.4±4.9%	38.2±2.6%	38.8±1.2%	42.2±2.4%	40.9±1.5%	85.2±0.8%	91.8±0.2%
Generalization gap		40.9%	27.2%	31.6%	15.8%	17.2%	3.8%	2.3%



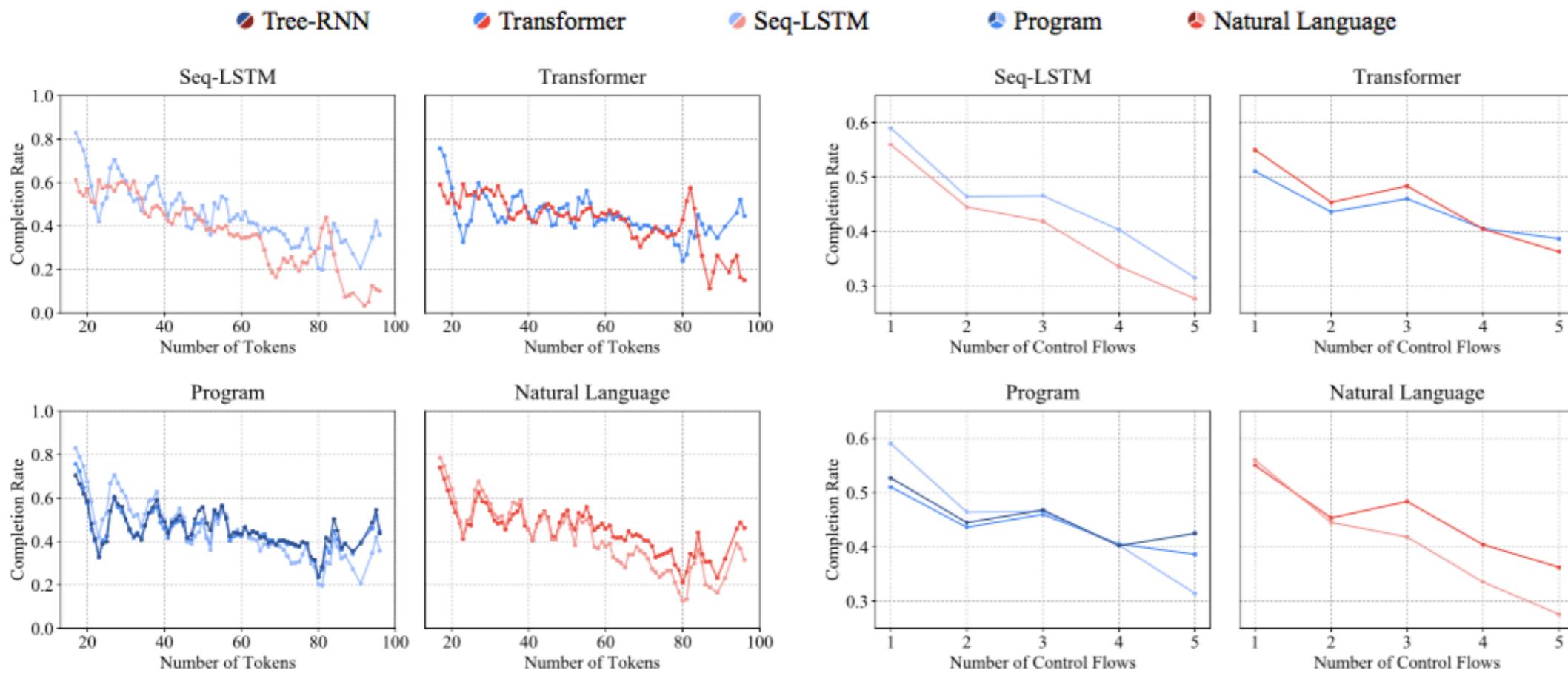
(a) Instruction Length

(b) Instruction Complexity

Quantitative Results

End-to-end < Ours (modular)

Instruction Method		Natural language descriptions		Programs				
		Seq-LSTM	Transformer	Seq-LSTM	Tree-RNN	Transformer	Ours (concat)	Ours
Dataset	test	54.9±1.8%	52.5±2.6%	56.7±1.9%	50.1±1.2%	49.4±1.6%	88.6±0.8%	94.0±0.5%
	test-complex	32.4±4.9%	38.2±2.6%	38.8±1.2%	42.2±2.4%	40.9±1.5%	85.2±0.8%	91.8±0.2%
Generalization gap		40.9%	27.2%	31.6%	15.8%	17.2%	3.8%	2.3%



(a) Instruction Length

(b) Instruction Complexity

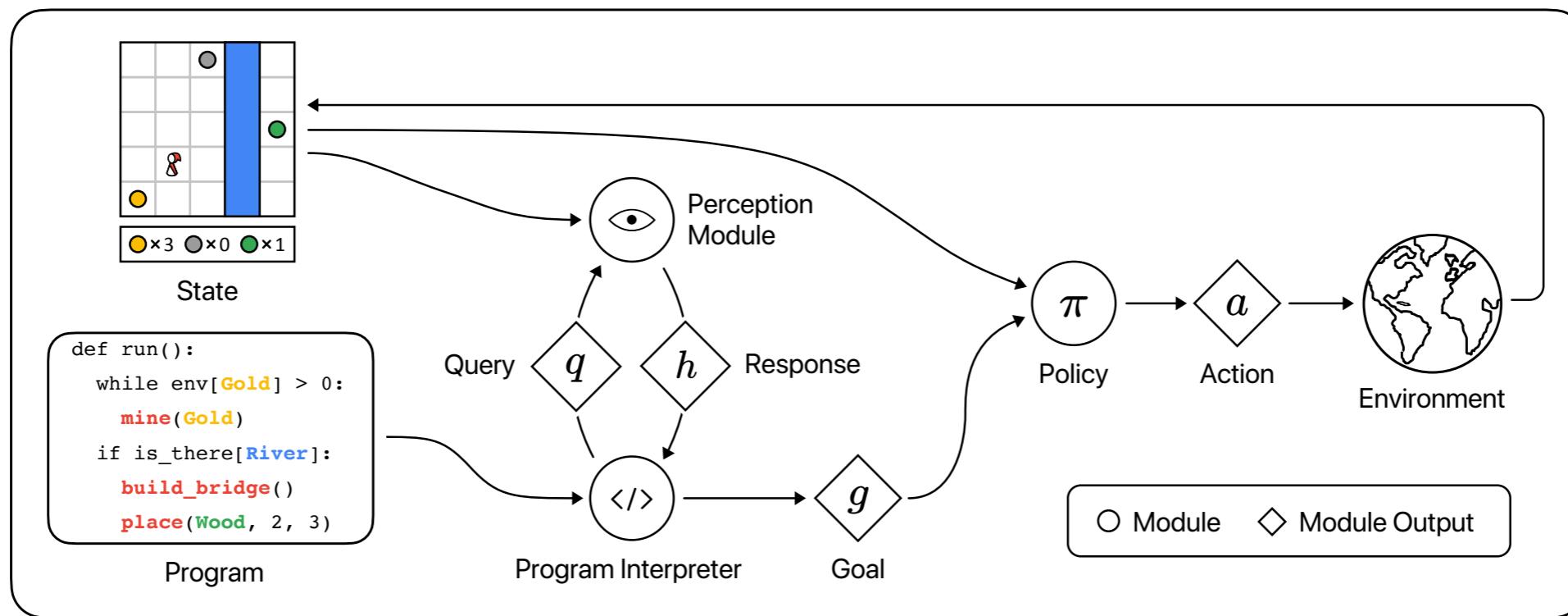
Takeaway

- Specific tasks using programs

Program

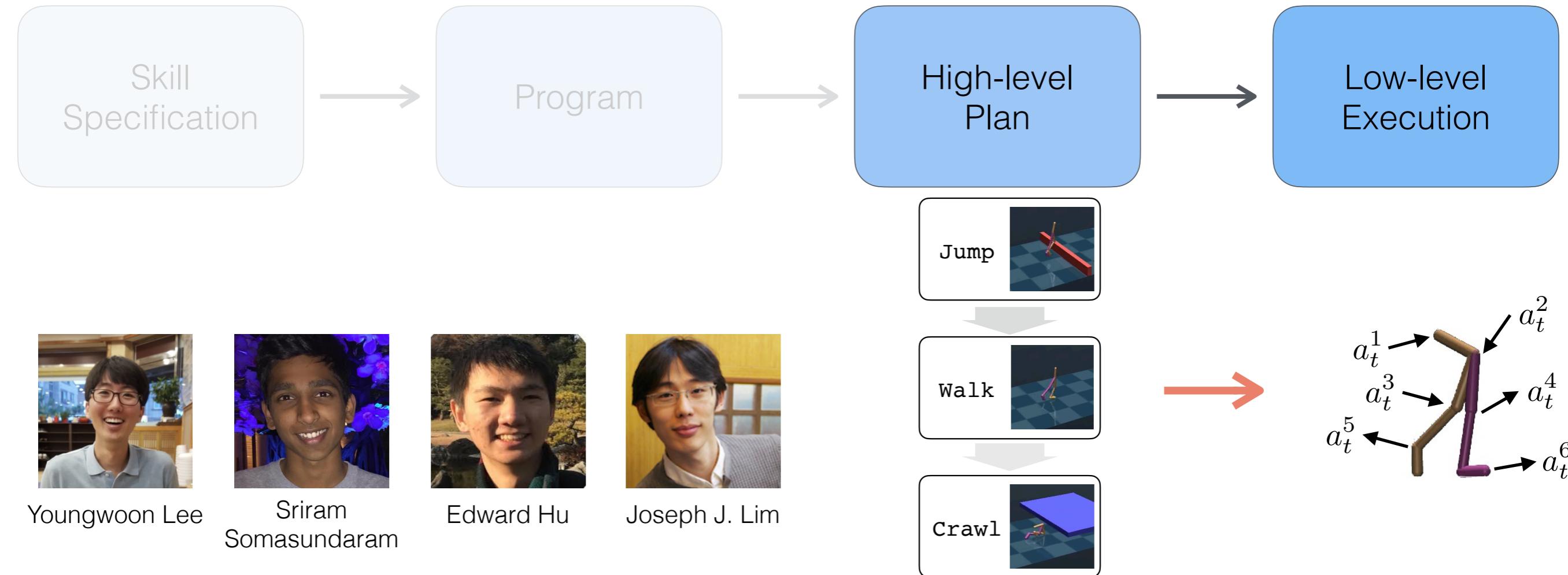
```
def run():
    if is_there[River]:
        mine(Wood)
        build_bridge()
        if agent[Iron]<3:
            mine(Iron)
            place(Iron, 1, 1)
        else:
            goto(4, 2)
    while env[Gold]>0:
        mine(Gold)
```

- Leverage the structure of programs with a modular framework

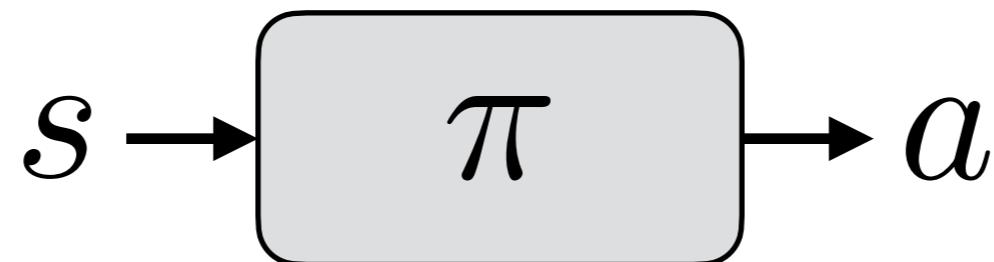


Composing Complex Skills by Learning Transition Policies

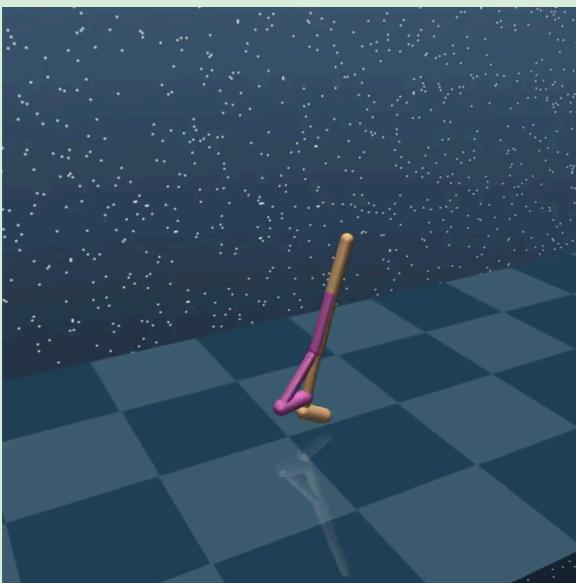
ICLR 2019



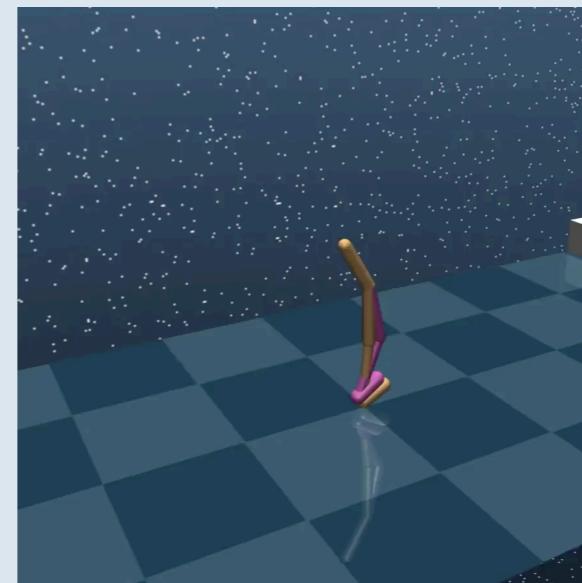
Learned Skills



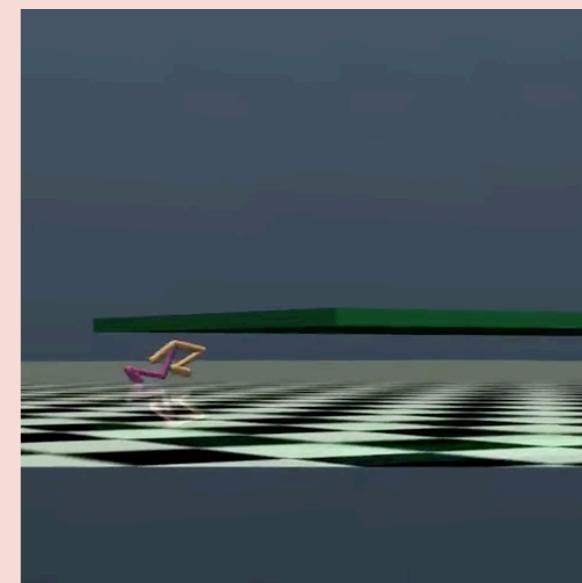
Walk



Jump



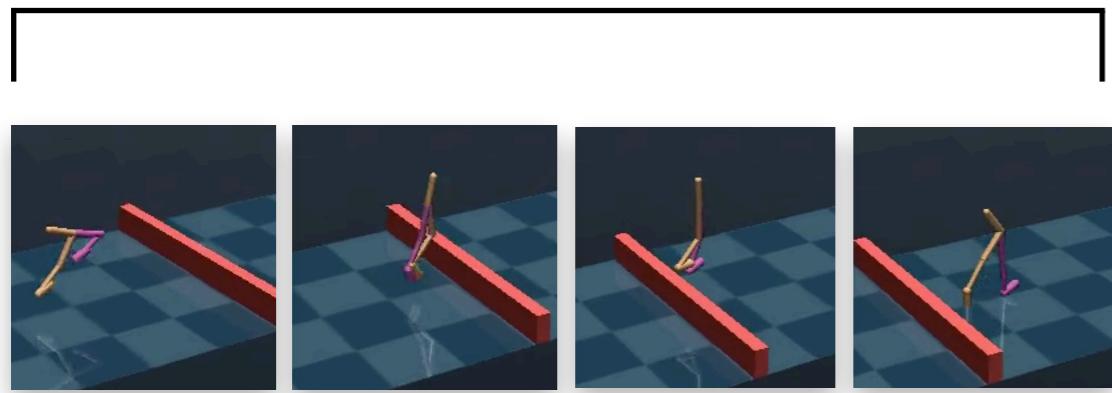
Crawl



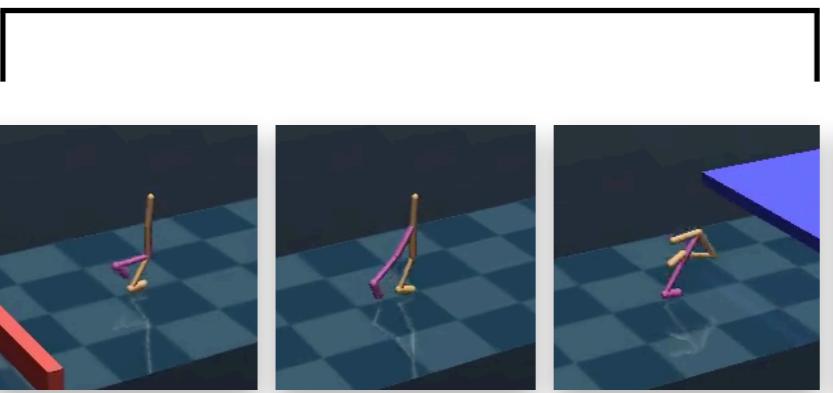
Compose Complex Skills

High-level plan

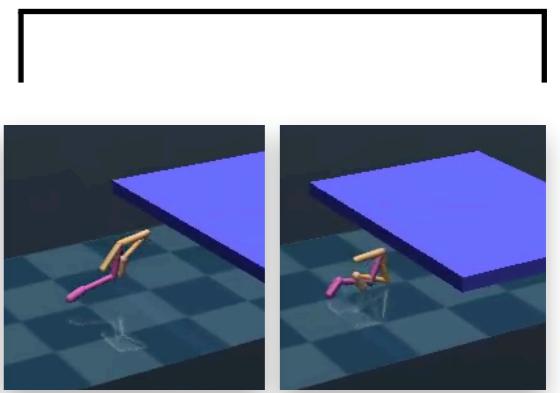
Jump



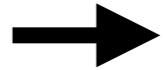
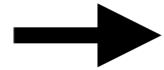
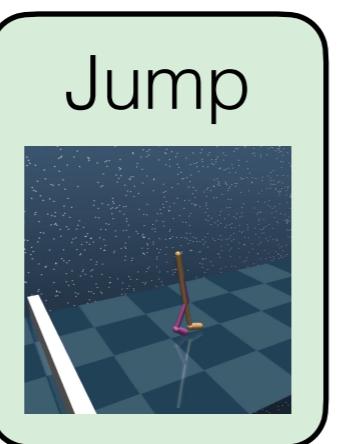
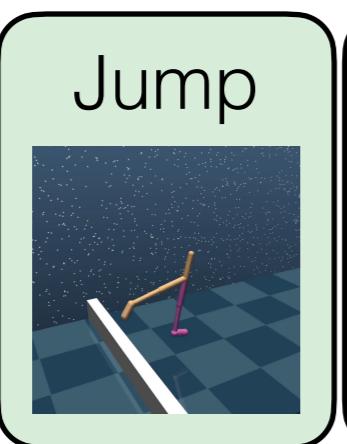
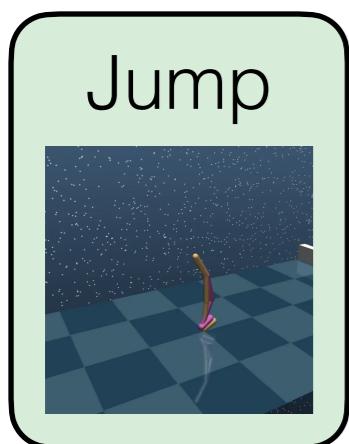
Walk



Crawl



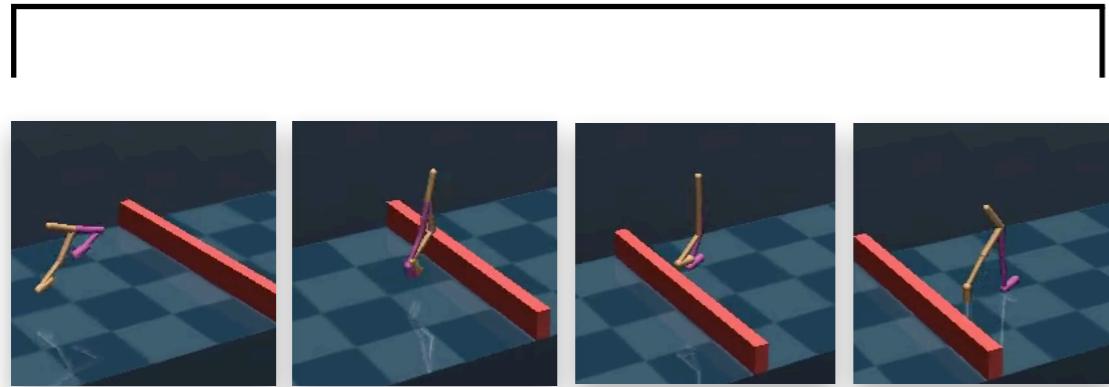
Sequentially execute corresponding policies



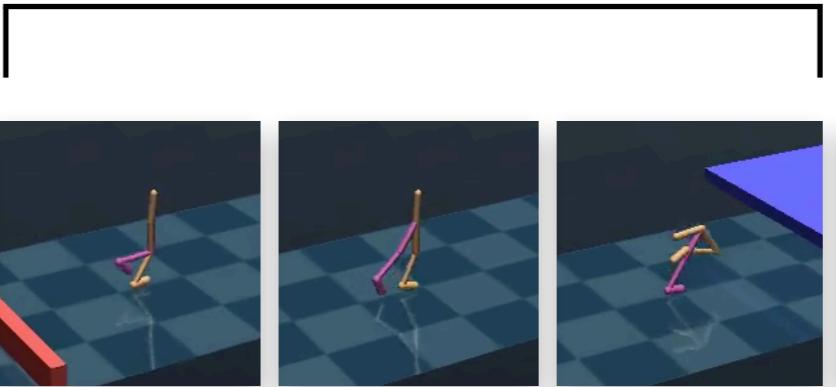
Compose Complex Skills

High-level plan

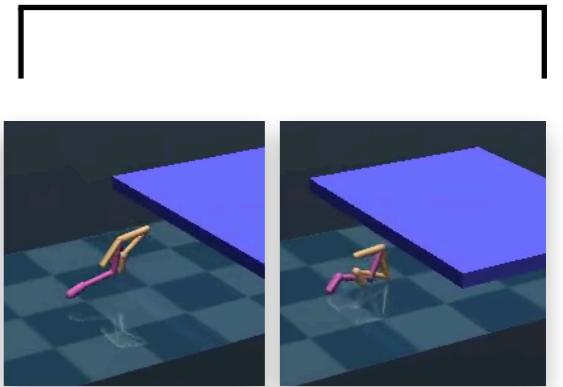
Jump



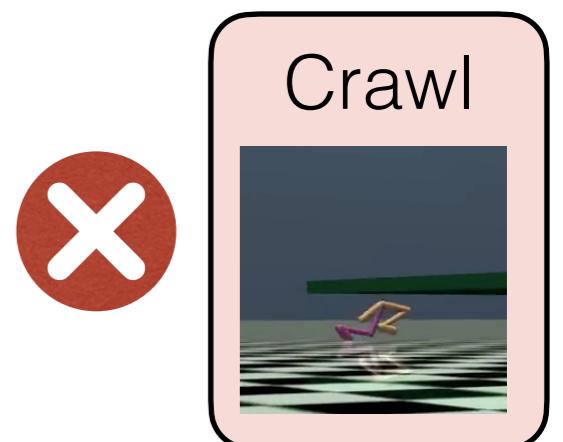
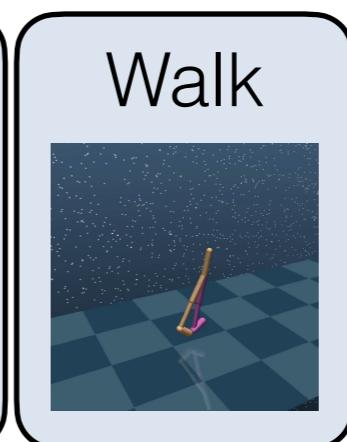
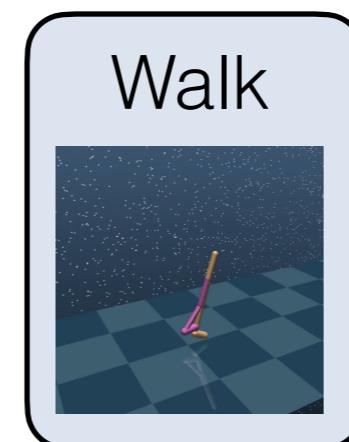
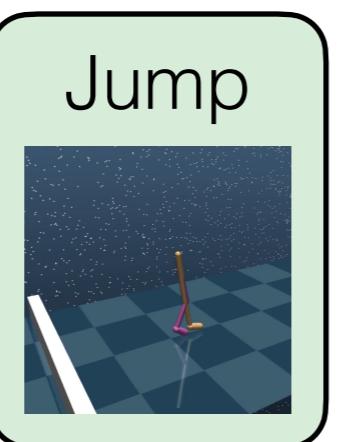
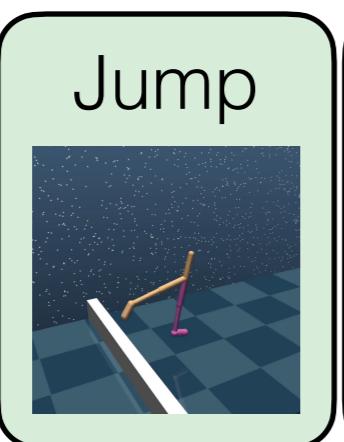
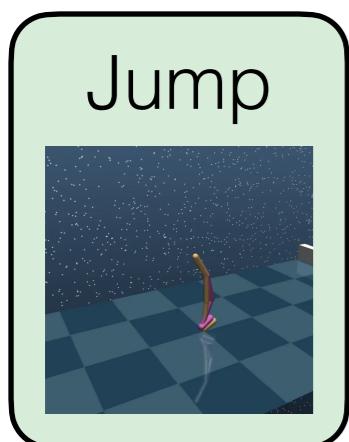
Walk



Crawl



Sequentially execute corresponding policies

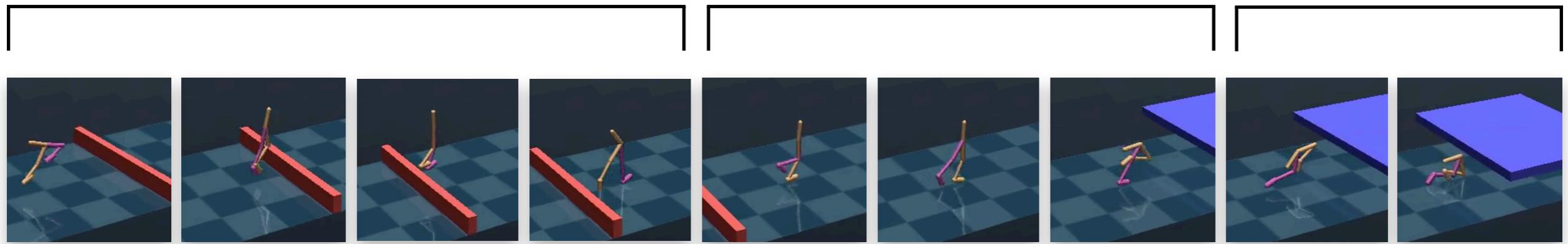


An end state of a previous policy might not be a good initial state of the following policy

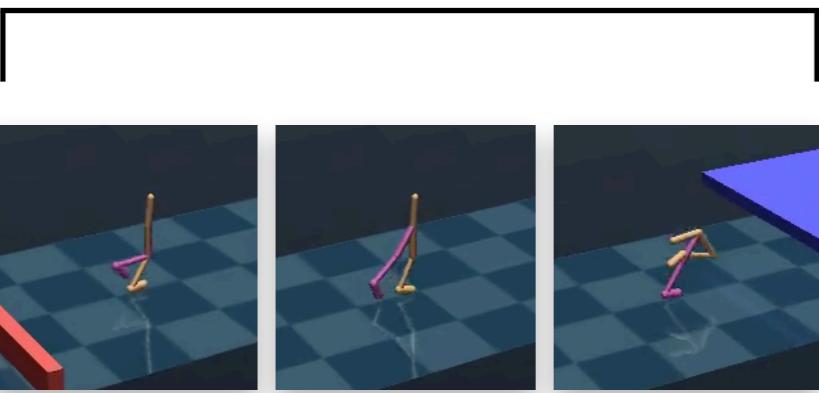
Compose Complex Skills

High-level plan

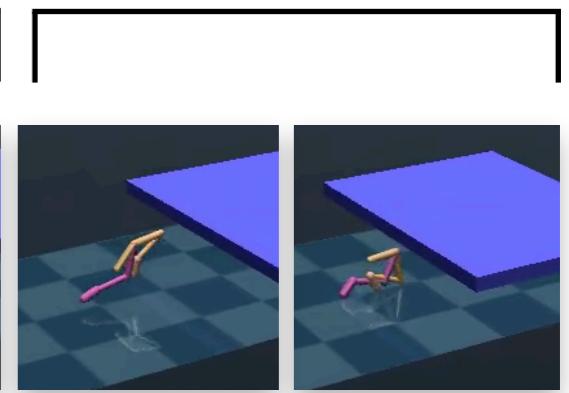
Jump



Walk



Crawl

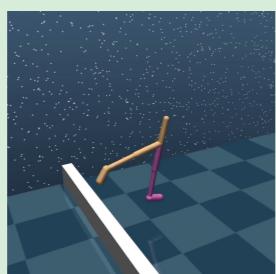


Sequentially execute corresponding policies

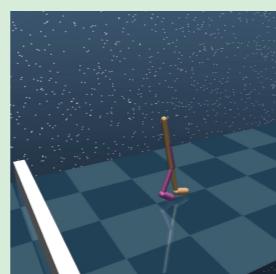
Jump



Jump

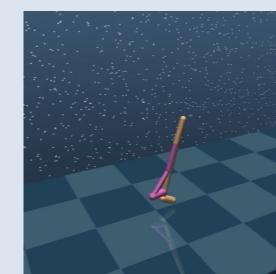


Jump

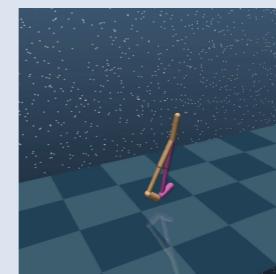


Trans
 π

Walk



Walk



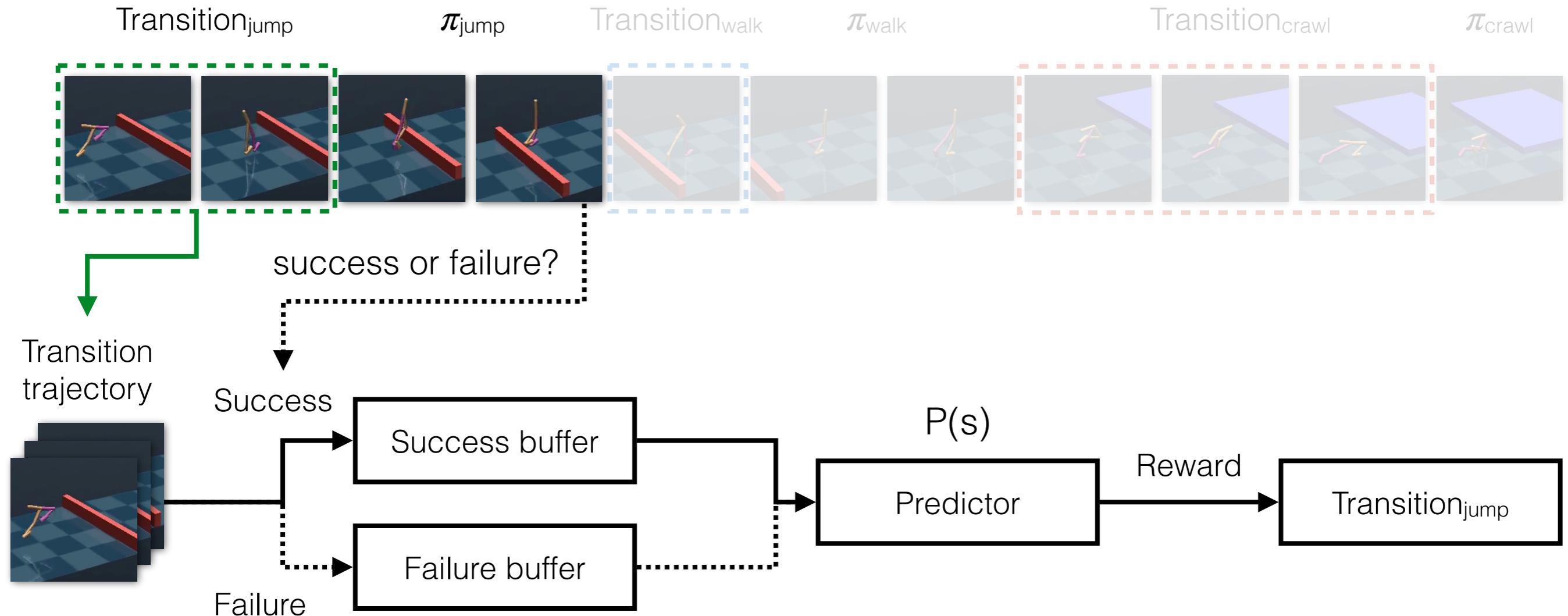
Trans
 π

Crawl



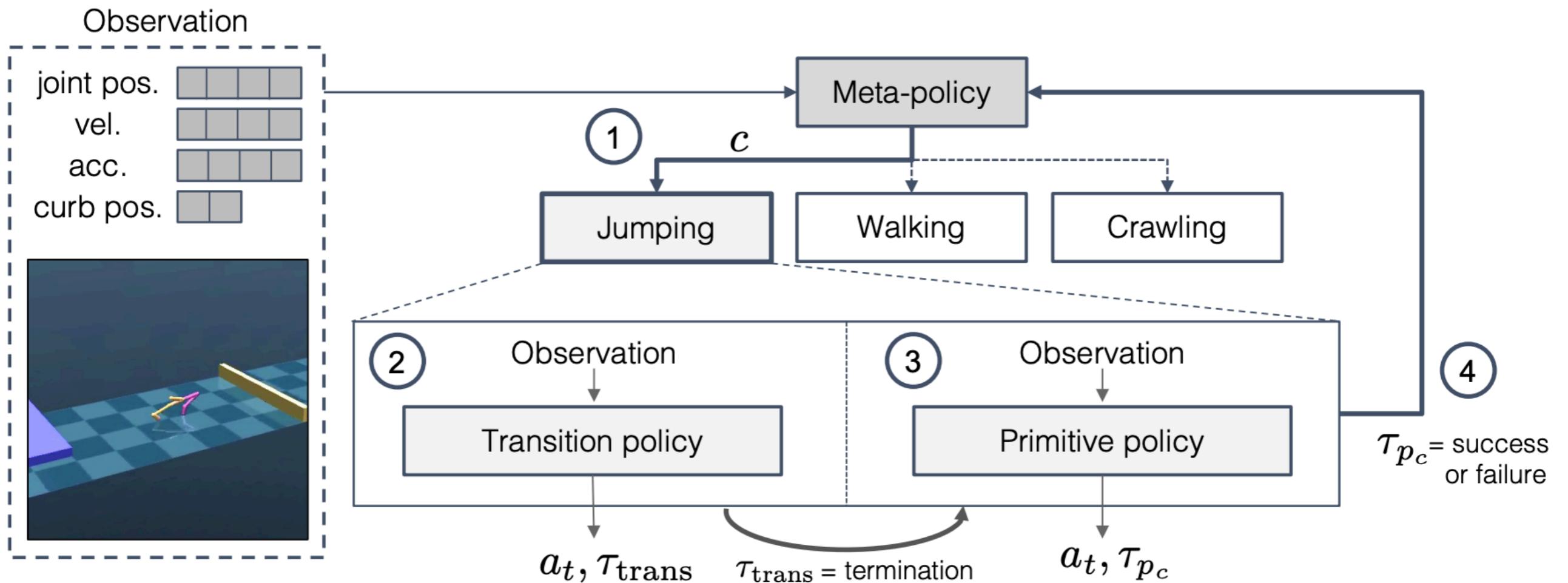
Transition policies

Learning Transition Policies



- Predictor learns to judge if a state is good for executing the next policy
- Transition policy learns from the predicted rewards

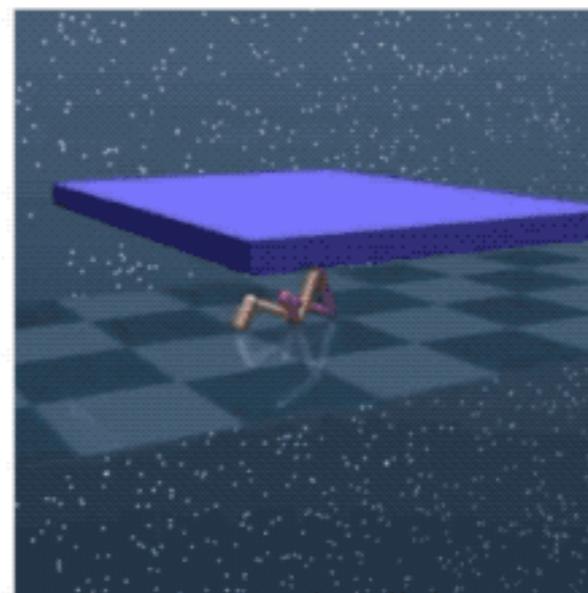
Modular Framework



Qualitative Results

Locomotion

Crawl



Transition

Walk

Walk Forward

Transition

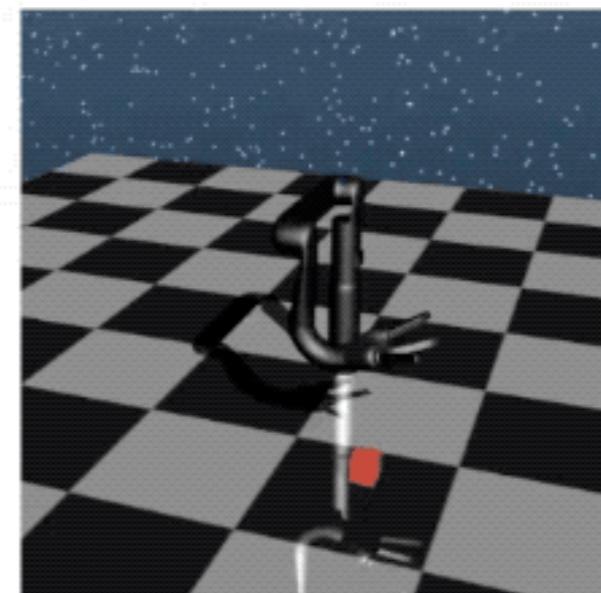
Walk Backward

Manipulation

Pick

Transition

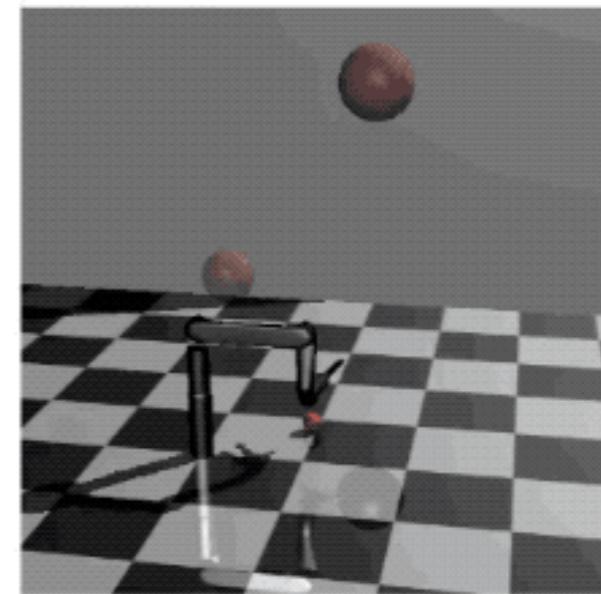
Pick



Toss

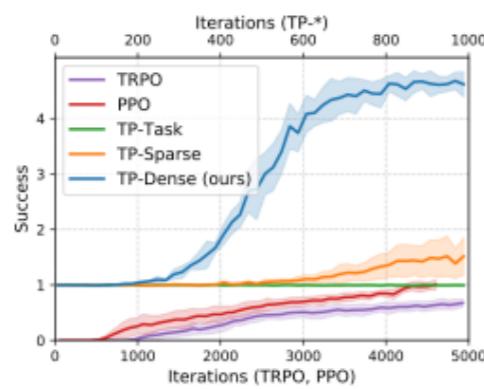
Transition

Hit

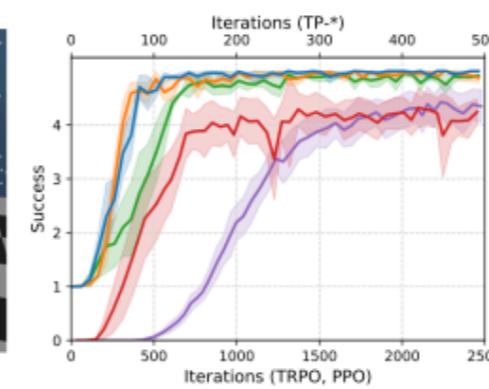


Quantitative Results

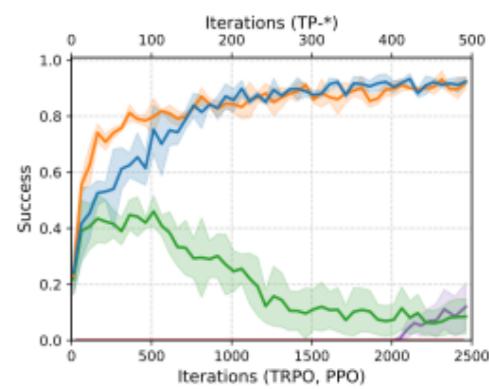
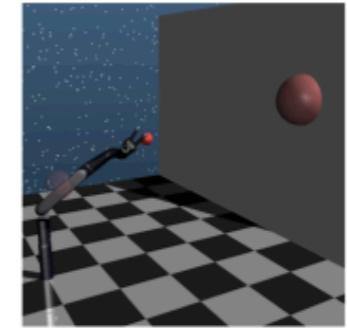
Sample Efficiency



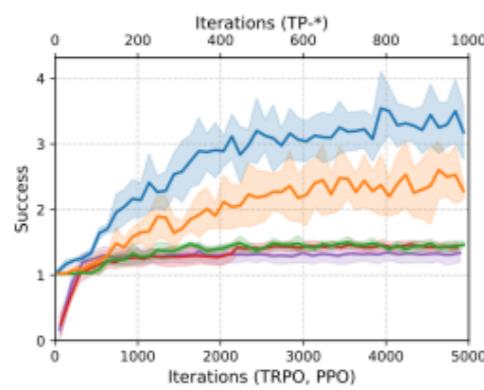
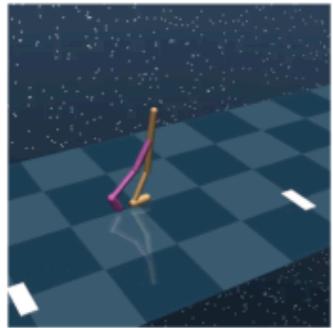
(a) Repetitive picking up



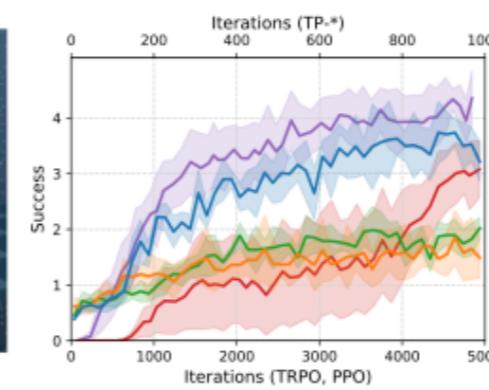
(b) Repetitive catching



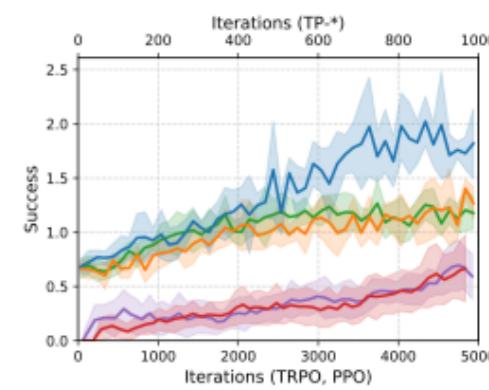
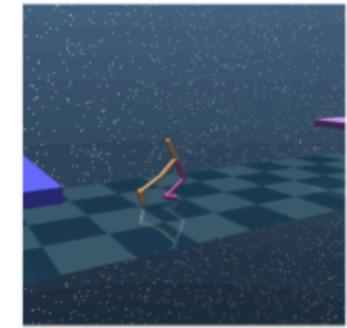
(c) Serve



(d) Patrol



(e) Hurdle

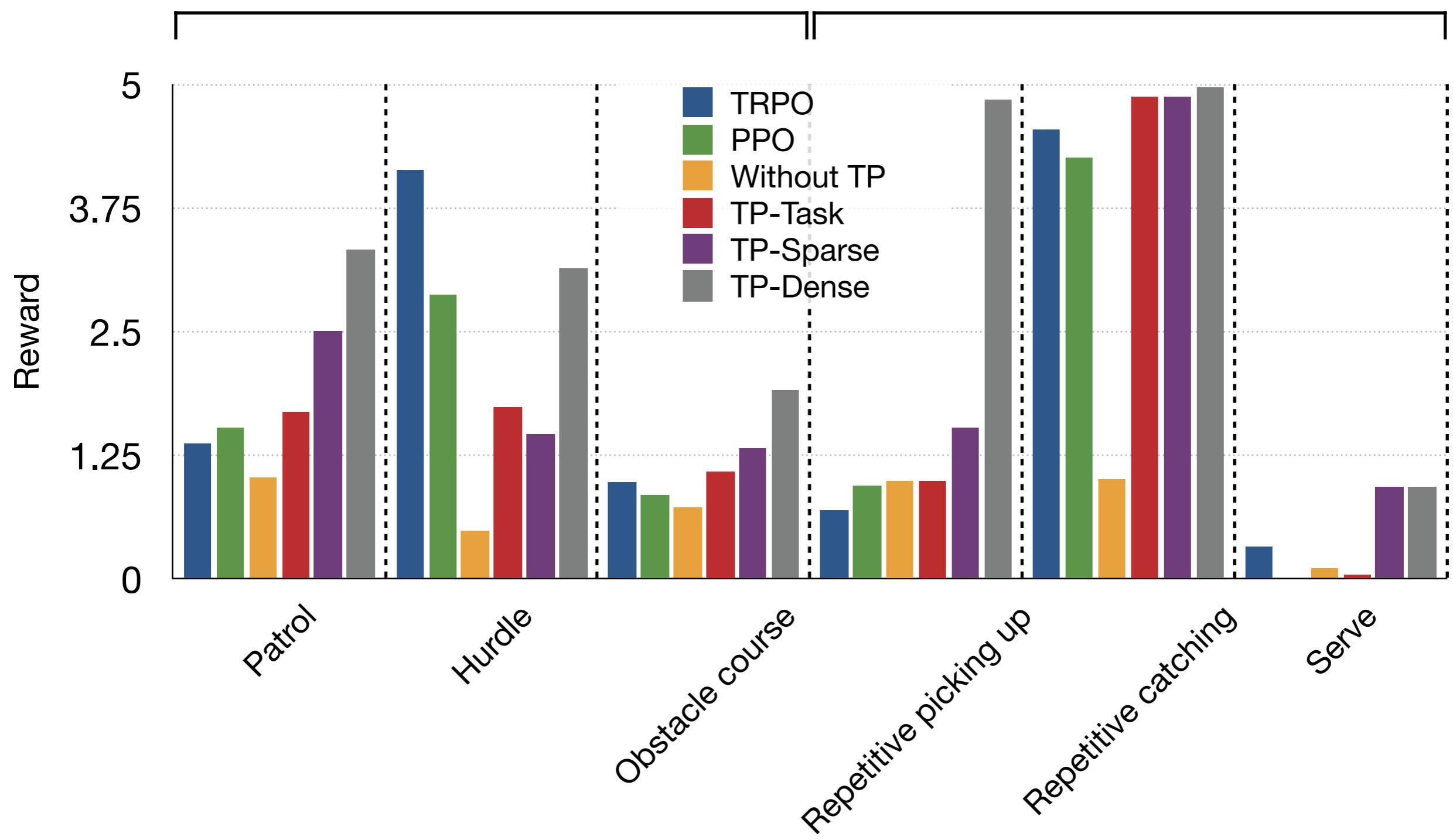


(f) Obstacle course

Quantitative Results

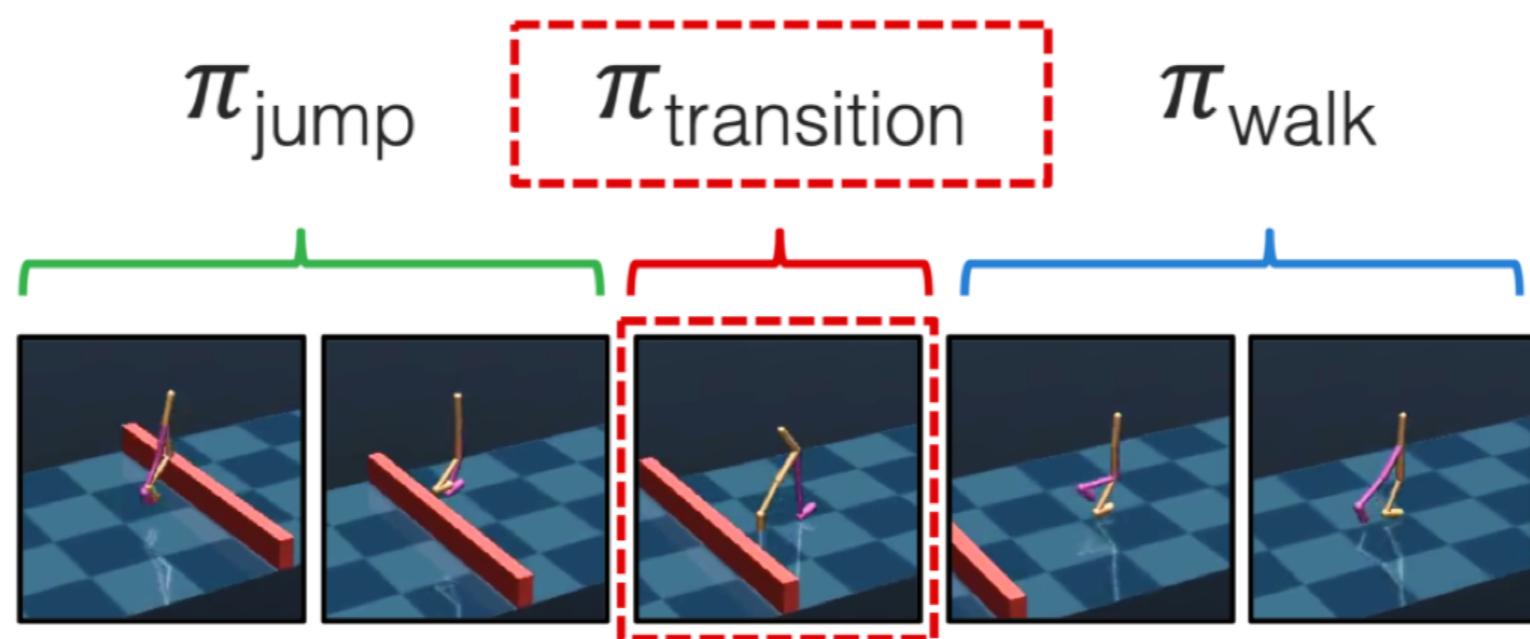
Locomotion

Manipulation

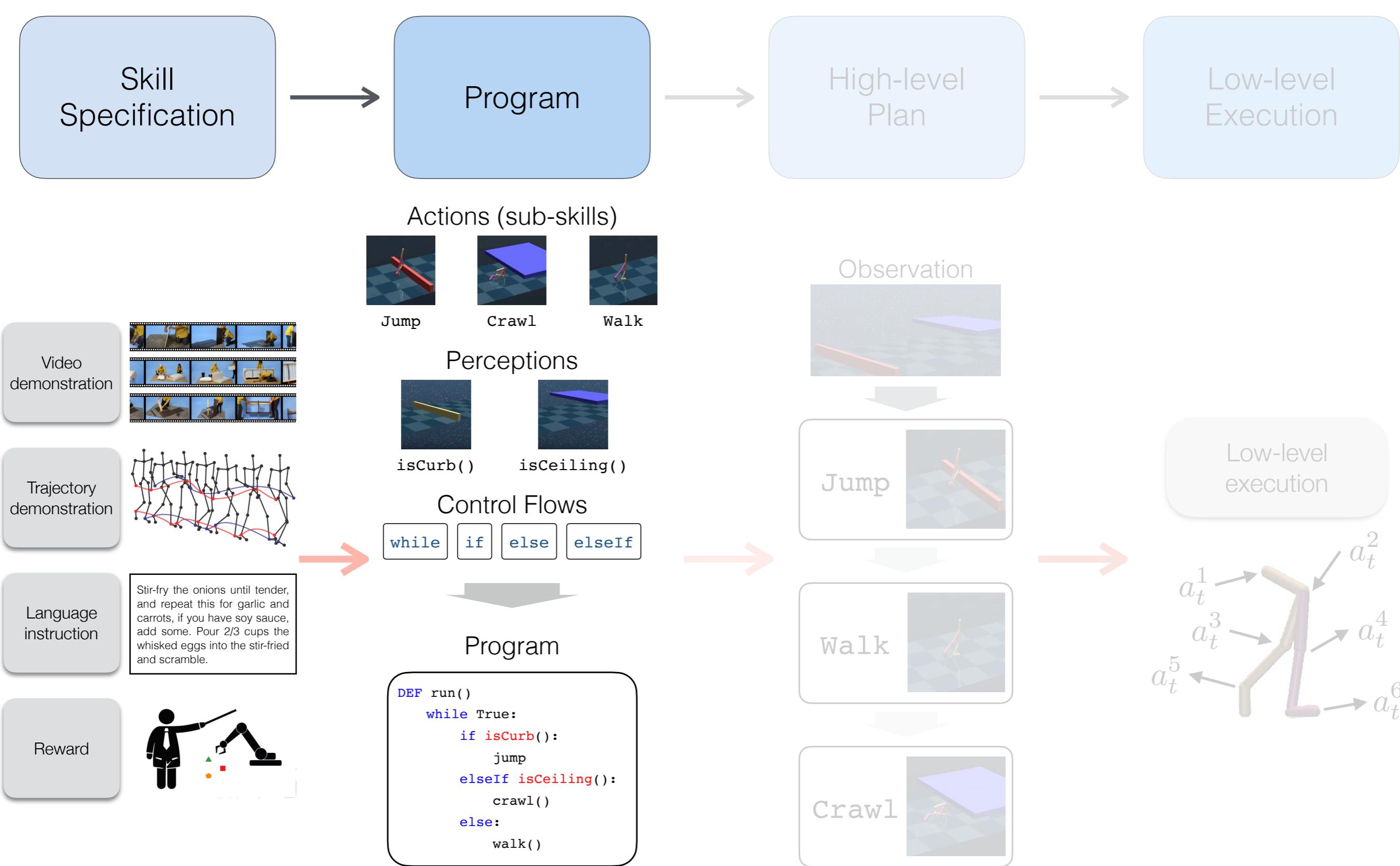


Takeaway

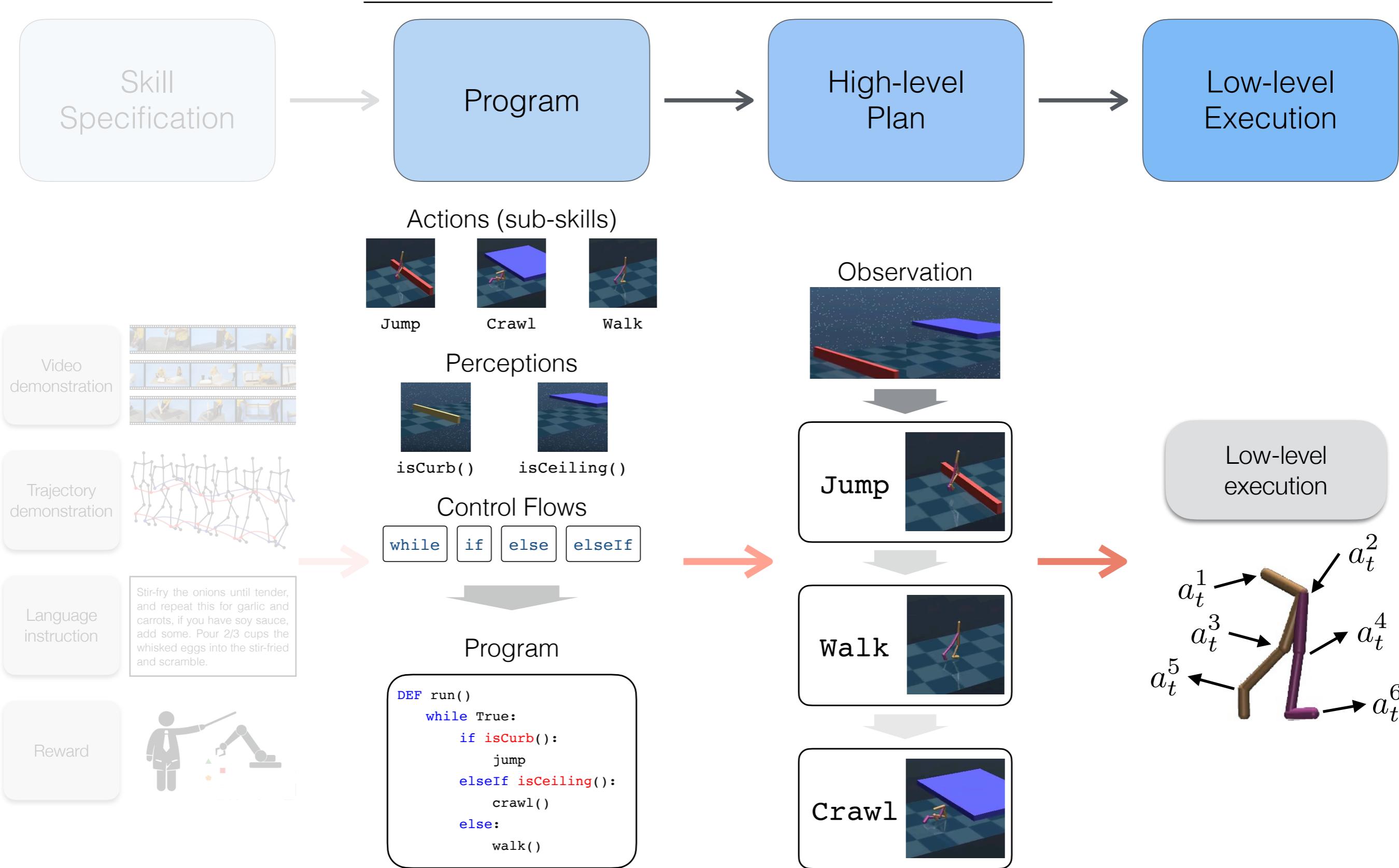
- Learning transition policies to smoothly compose learned skills



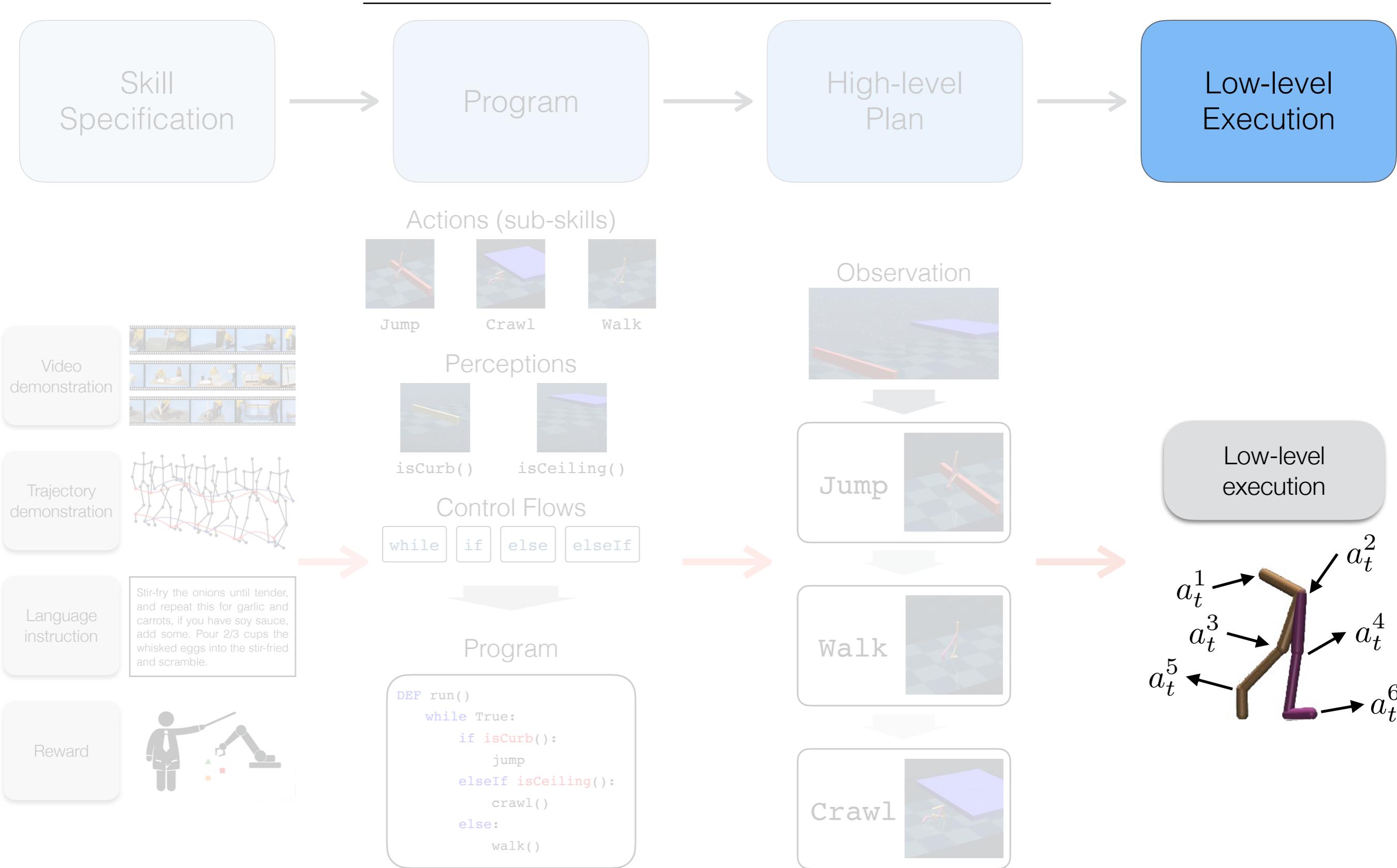
Program Inference



Task Execution

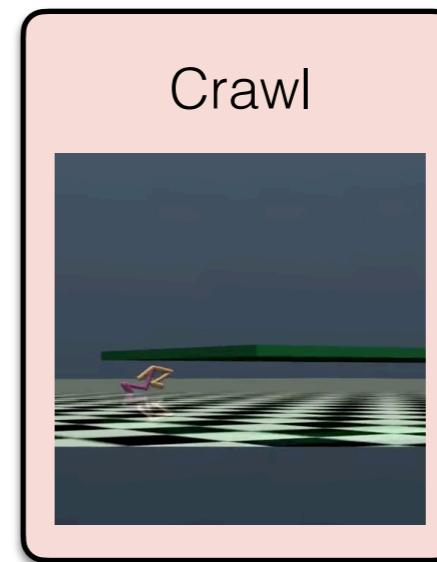
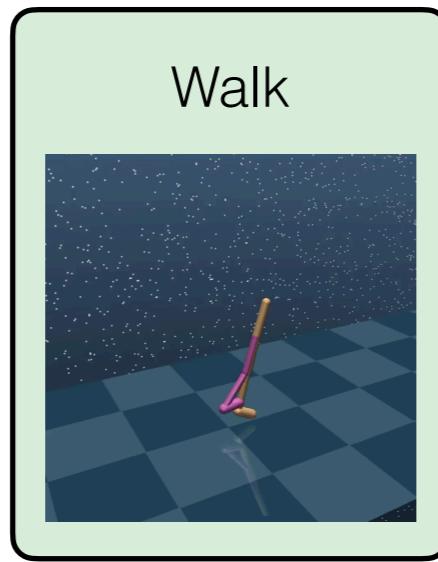


Primitive Skill Acquisition



Primitive Skill Acquisition

- Goal: acquire a diverse set of primitive skills efficiently



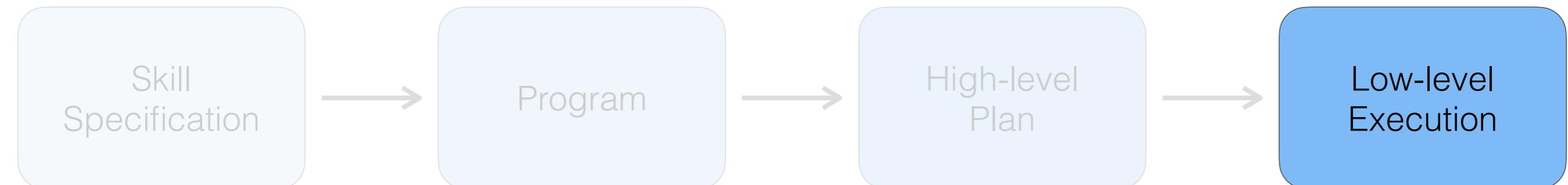
- Key directions

Meta-learning
Meta-RL

Learning from
experts

Toward Multimodal Model-Agnostic Meta-Learning

Meta-learning Workshop at NeurIPS 2018



Risto Vuorio

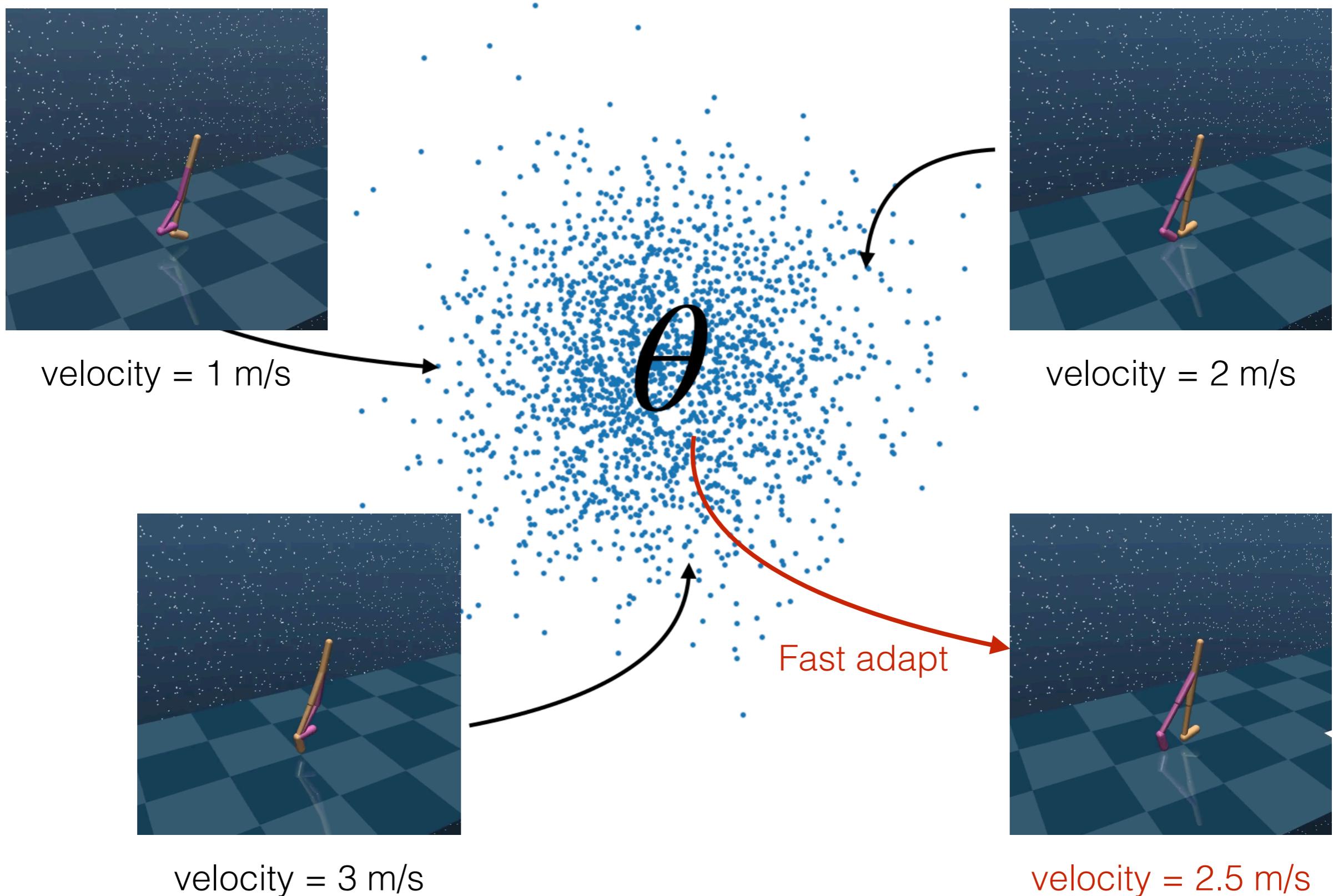


Hexiang Hu

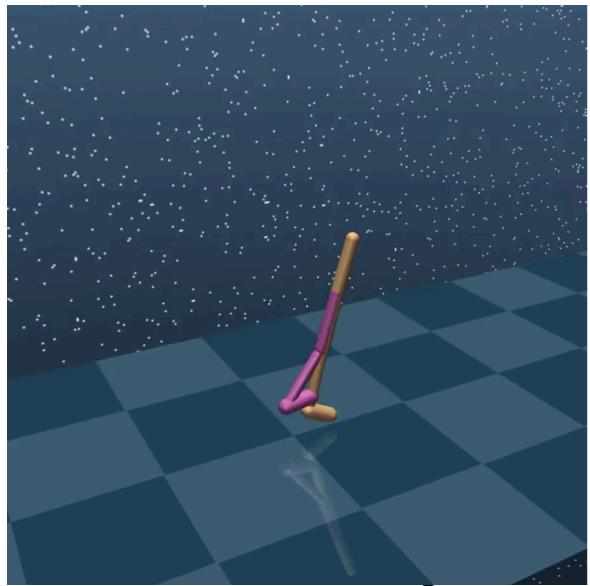


Joseph J. Lim

Model-Agnostic Meta-Learning (MAML)

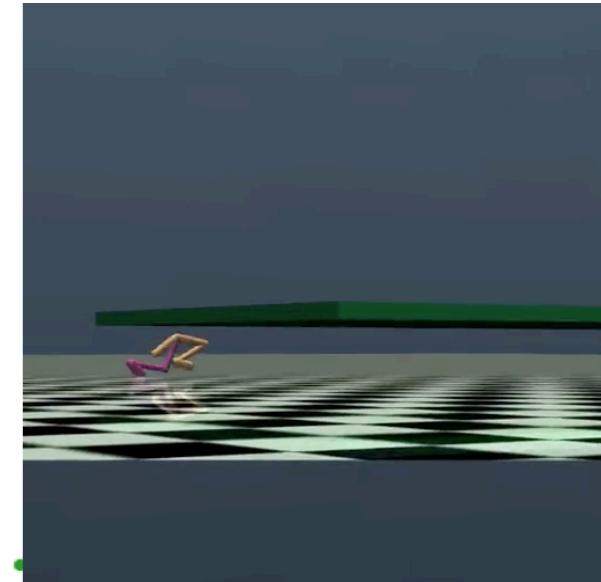


Walk

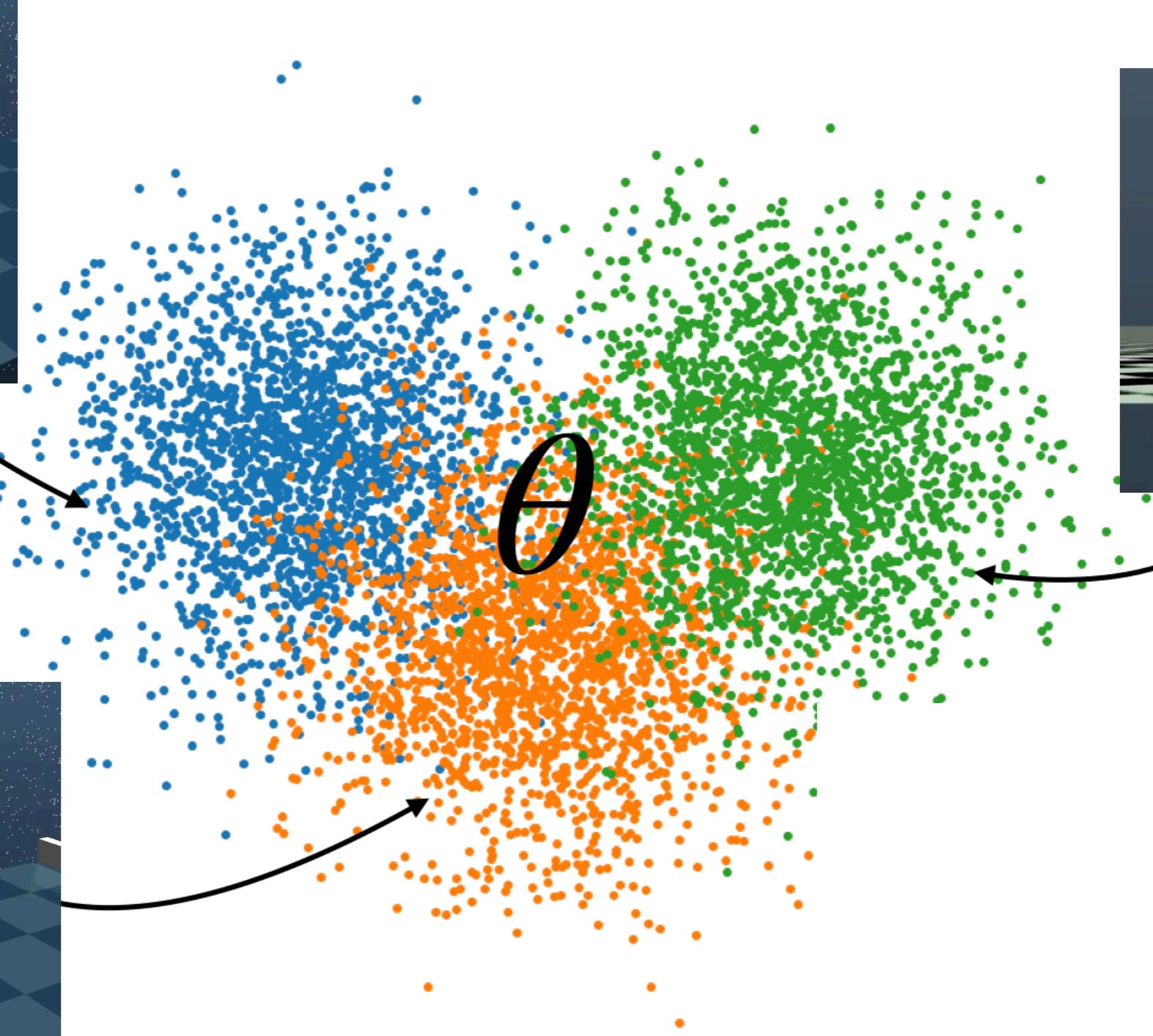
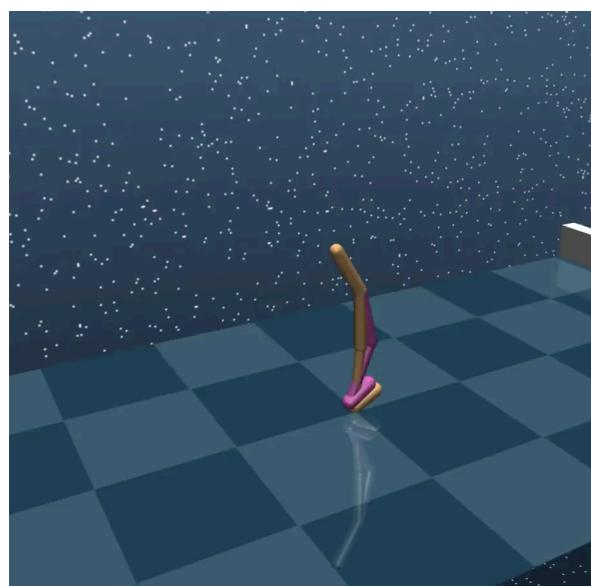


Multimodal Task Distribution

Crawl

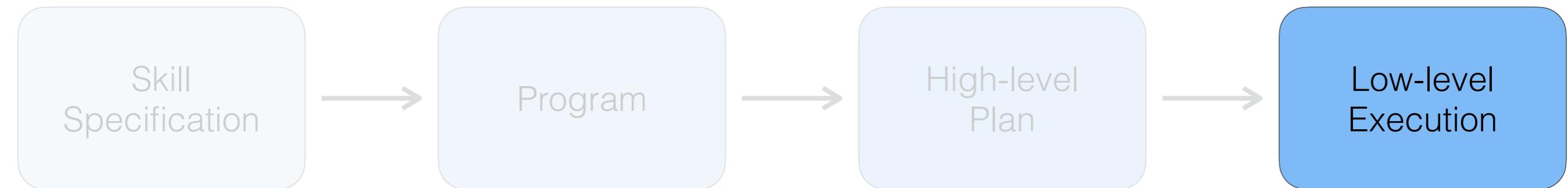


Jump



Multimodal Model-Agnostic Meta-Learning via Task-Aware Modulation

NeurIPS 2019 (Spotlight)



Risto Vuorio

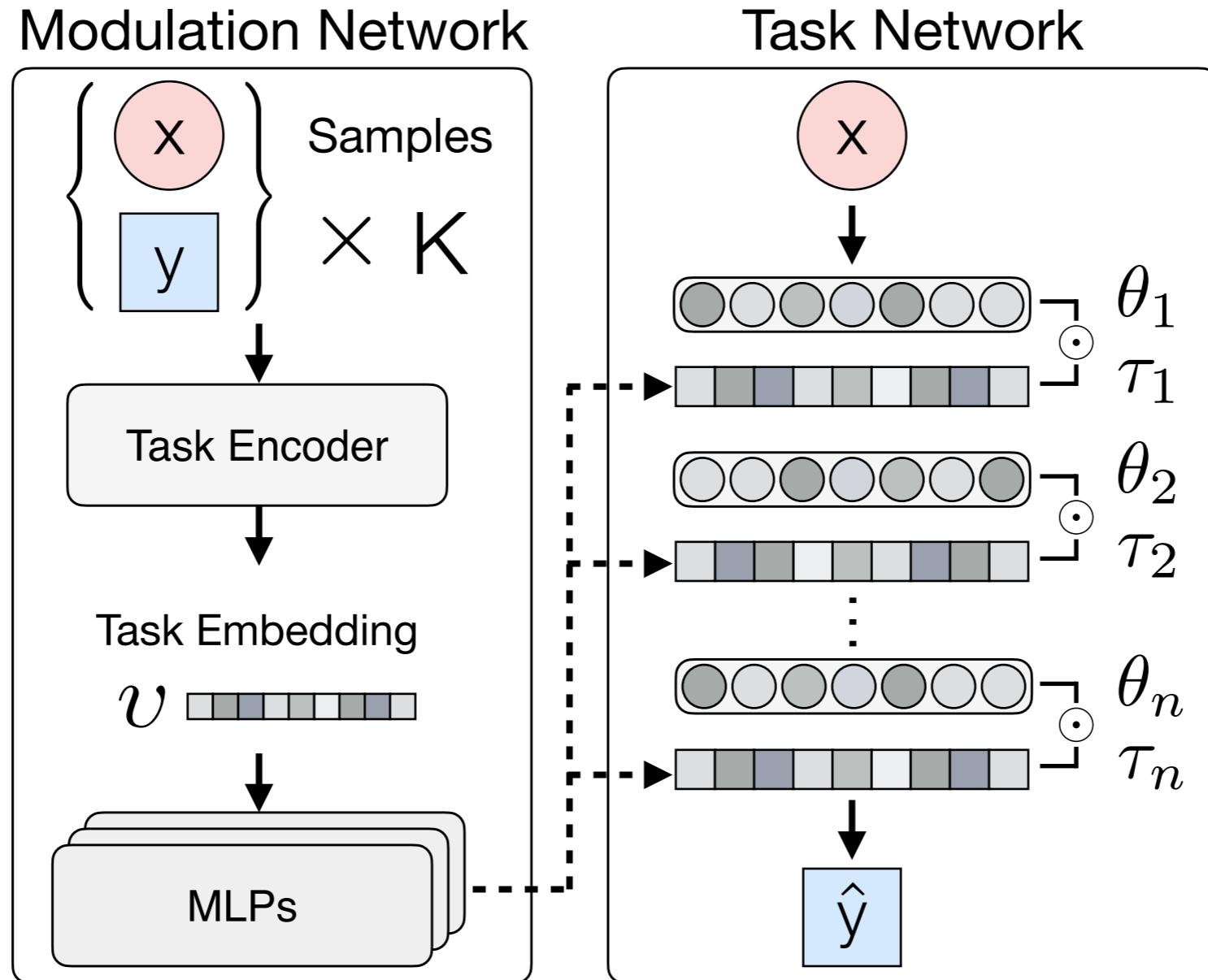


Hexiang Hu



Joseph J. Lim

Multimodal Model-Agnostic Meta-Learning (MMAML)



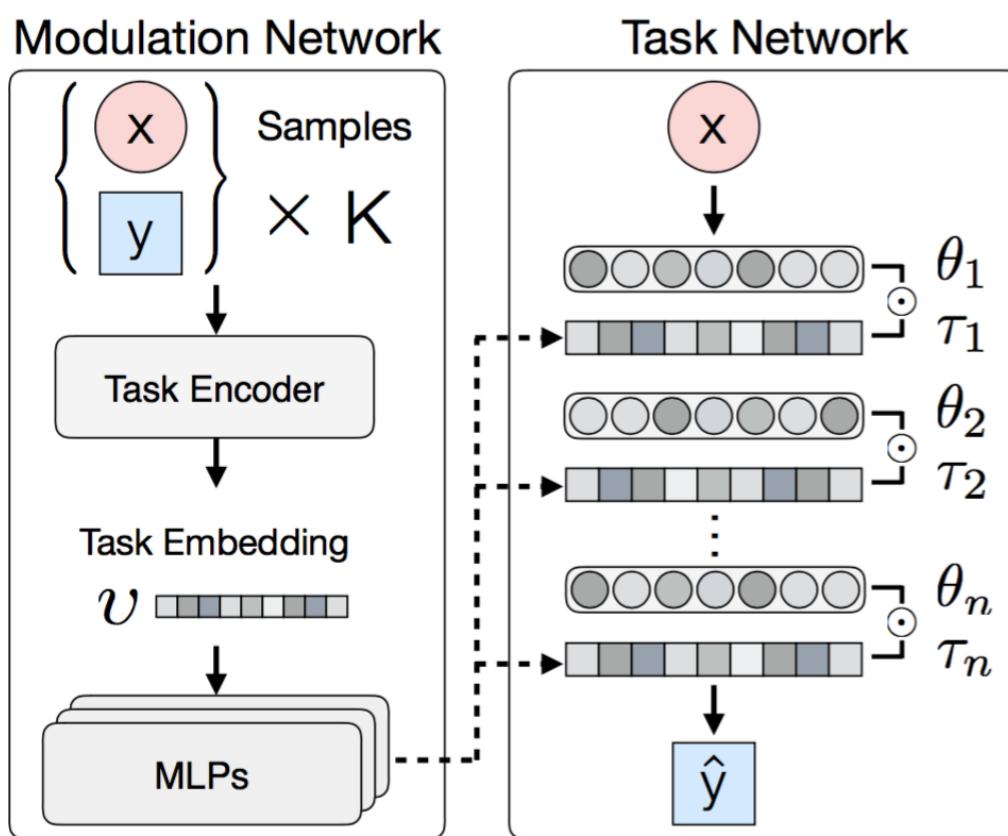
Training Algorithm

Outer loop

- Task Encoder: produce the task embedding ω_g
- MLPs: modulate the task network blocks ω_h

Inner loop

- Task network: fast adapt through gradient updates θ



Algorithm 1 MMAML META-TRAINING PROCEDURE.

```

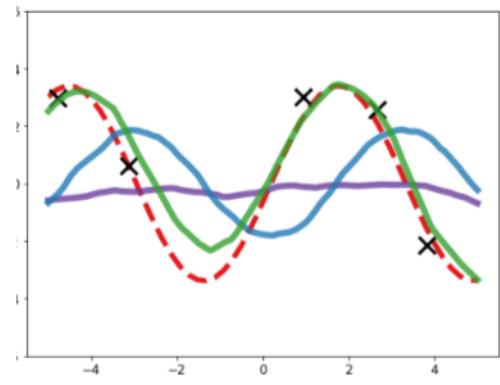
1: Input: Task distribution  $P(\mathcal{T})$ , Hyper-parameters  $\alpha$  and  $\beta$ 
2: Randomly initialize  $\theta$  and  $\omega$ .
3: while not DONE do
4:   Sample batches of tasks  $\mathcal{T}_j \sim P(\mathcal{T})$ 
5:   for all j do
6:     Infer  $v = h(\{x, y\}_K; \omega_h)$  with K samples from  $\mathcal{D}_{\mathcal{T}_j}^{\text{train}}$ .
7:     Generate parameters  $\tau = \{g_i(v; \omega_g) \mid i = 1, \dots, N\}$  to modulate each block of the task network  $f$ .
8:     Evaluate  $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_j}(f(x; \theta, \tau); \mathcal{D}_{\mathcal{T}_j}^{\text{train}})$  w.r.t the K samples
9:     Compute adapted parameter with gradient descent:
10:     $\theta'_{\mathcal{T}_j} = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_j}(f(x; \theta, \tau); \mathcal{D}_{\mathcal{T}_j}^{\text{train}})$ 
11:   end for
12:   Update  $\theta$  with  $\beta \nabla_{\theta} \sum_{T_j \sim P(\mathcal{T})} \mathcal{L}_{\mathcal{T}_j}(f(x; \theta', \tau); \mathcal{D}_{\mathcal{T}_j}^{\text{val}})$ 
13:   Update  $\omega_g$  with  $\beta \nabla_{\omega_g} \sum_{T_j \sim P(\mathcal{T})} \mathcal{L}_{\mathcal{T}_j}(f(x; \theta', \tau); \mathcal{D}_{\mathcal{T}_j}^{\text{val}})$ 
14: end while

```

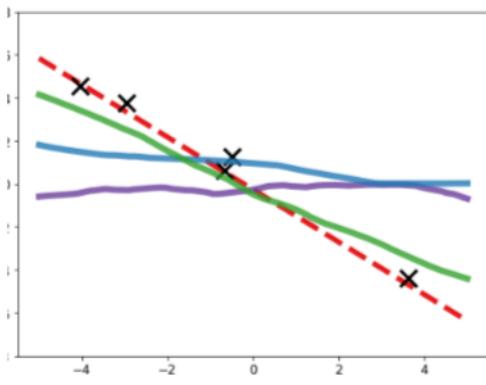
Regression



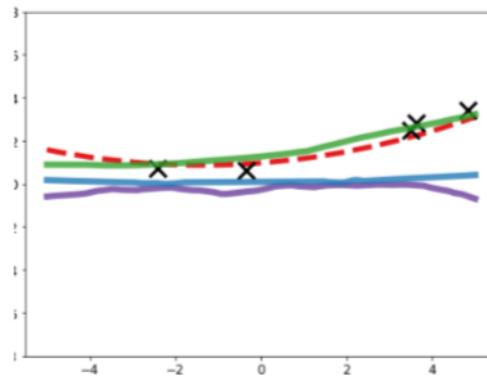
Sinusoidal



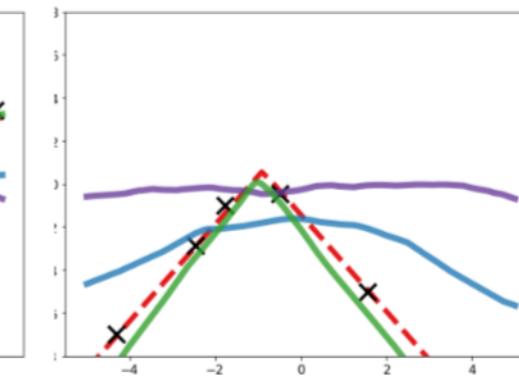
Linear



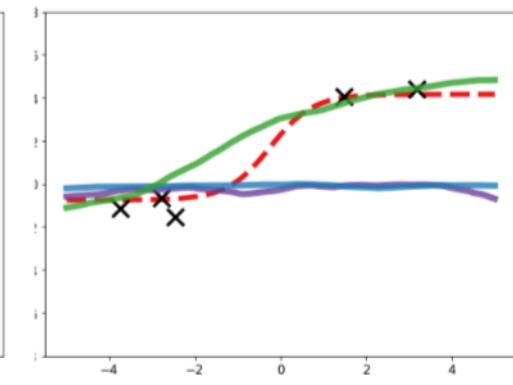
Quadratic



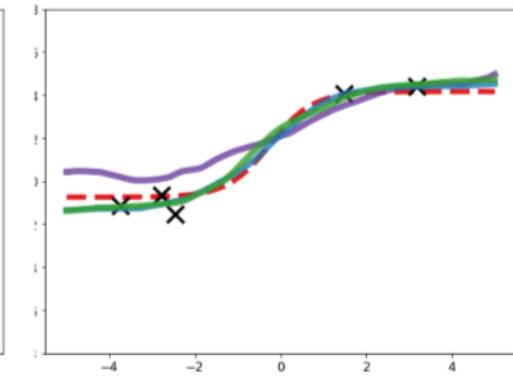
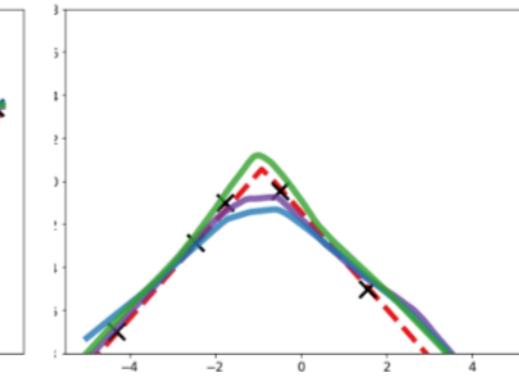
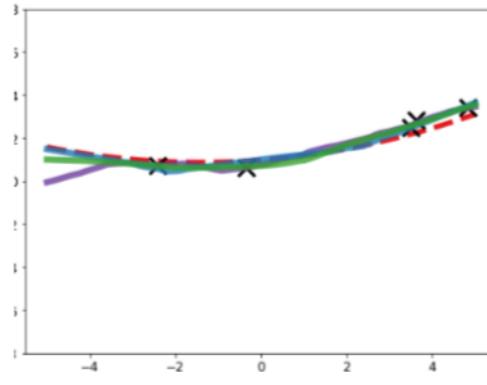
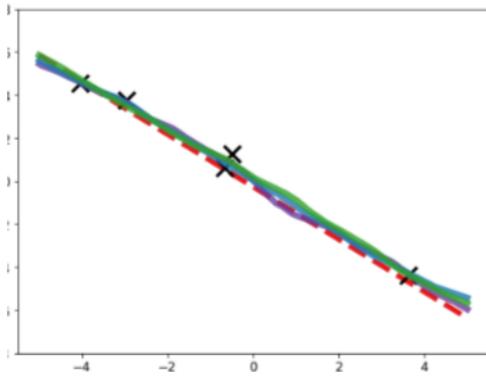
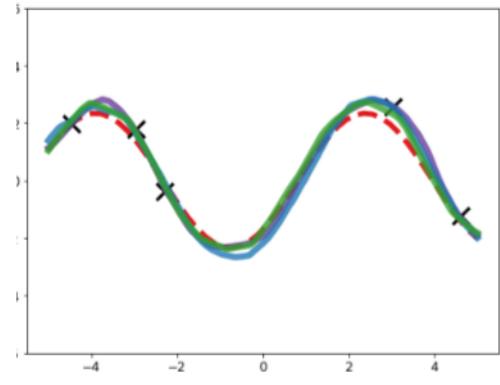
Transformed ℓ_1 Norm



Tanh



(a) MMAML post modulation vs. other prior models

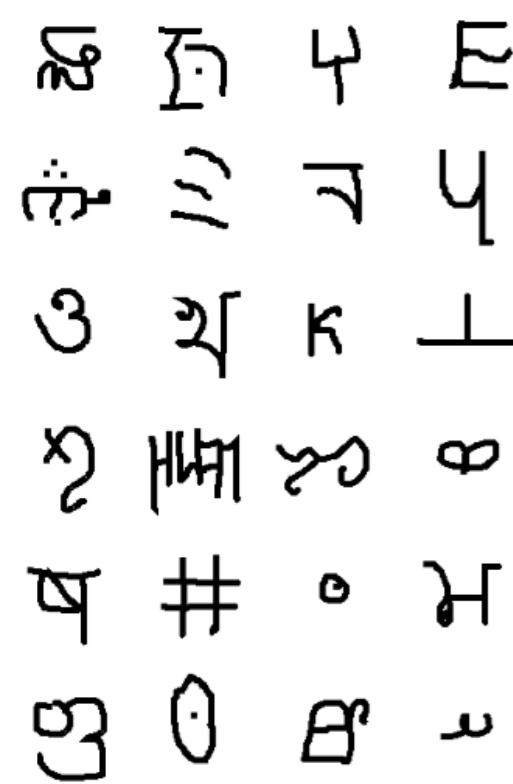


(b) MMAML post adaptation vs. other posterior models

Method	2 Modes		3 Modes		5 Modes	
	Post Modulation	Post Adaptation	Post Modulation	Post Adaptation	Post Modulation	Post Adaptation
MAML [1]	-	1.085	-	1.231	-	1.668
Multi-MAML	-	0.433	-	0.713	-	1.082
LSTM Learner	0.362	-	0.548	-	0.898	-
Ours: MMAML (Softmax)	1.548	0.361	2.213	0.444	2.421	0.939
Ours: MMAML (FiLM)	2.421	0.336	1.923	0.444	2.166	0.868

Image Classification

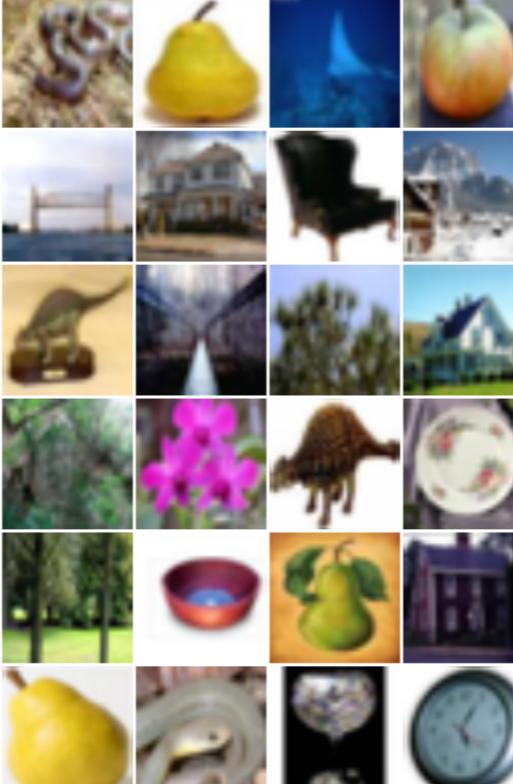
(a) Omniglot



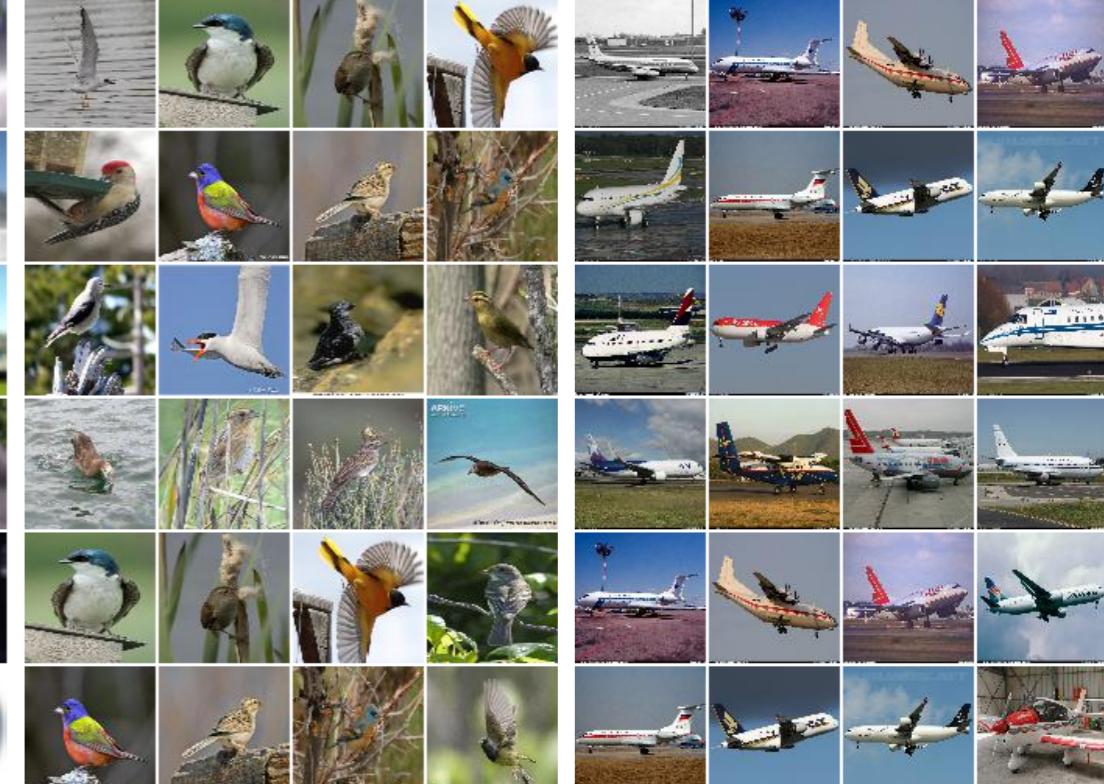
(b) Mini-ImageNet



(c) FC100



(d) CUB



(e) Aircraft



Method & Setup

2 Modes

3 Modes

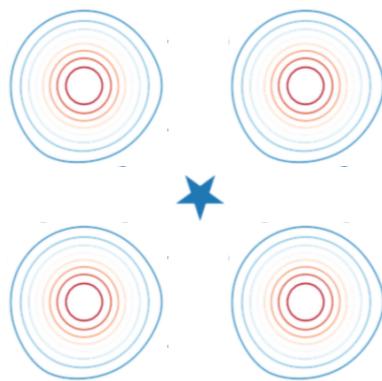
5 Modes

Way Shot	5-way 1-shot		20-way 1-shot		5-way 1-shot		20-way 1-shot		5-way 1-shot		20-way 3-shot	
MAML [1]	66.80%	77.79%	44.69%	54.55%	67.97%	28.22%	66.00%	70.87%	44.69%	56.57%		
Multi-MAML	66.85%	73.07%	53.15%	55.90%	62.20%	39.77%	71.94%	76.94%	58.66%	61.11%		
MMAML (ours)	69.93%	78.73%	47.80%	57.47%	70.15%	36.27%	72.02%	78.13%	60.13%	63.52%		

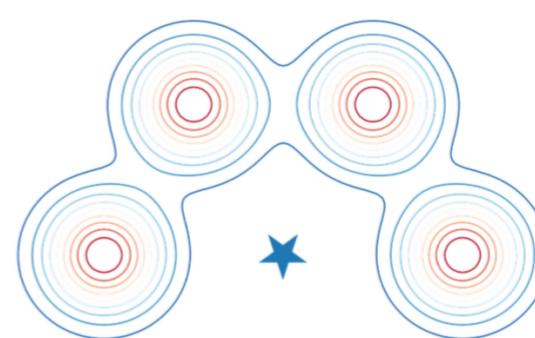
Reinforcement Learning

Goal modes

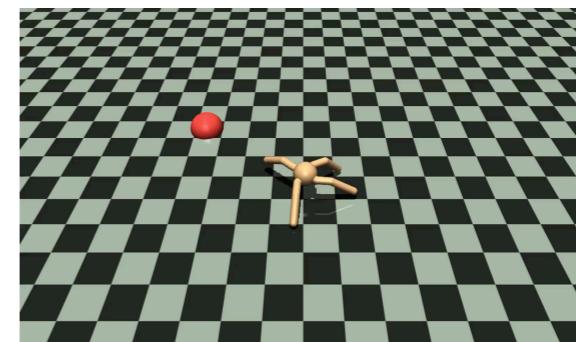
Reacher



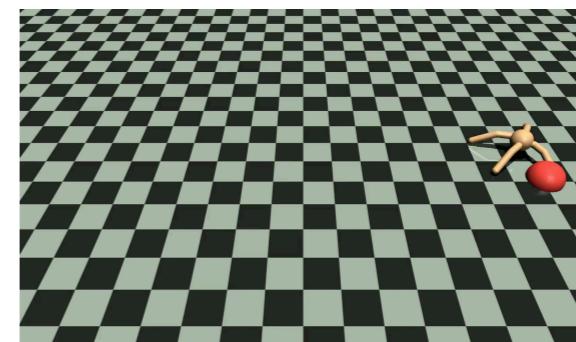
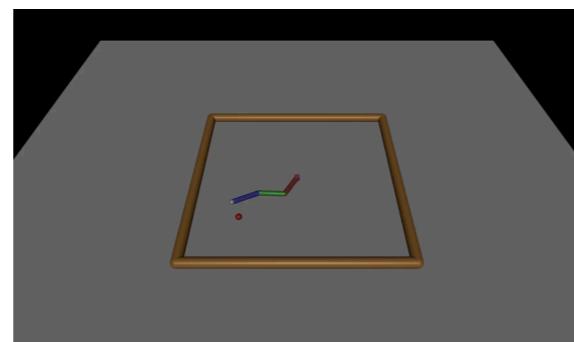
Ant



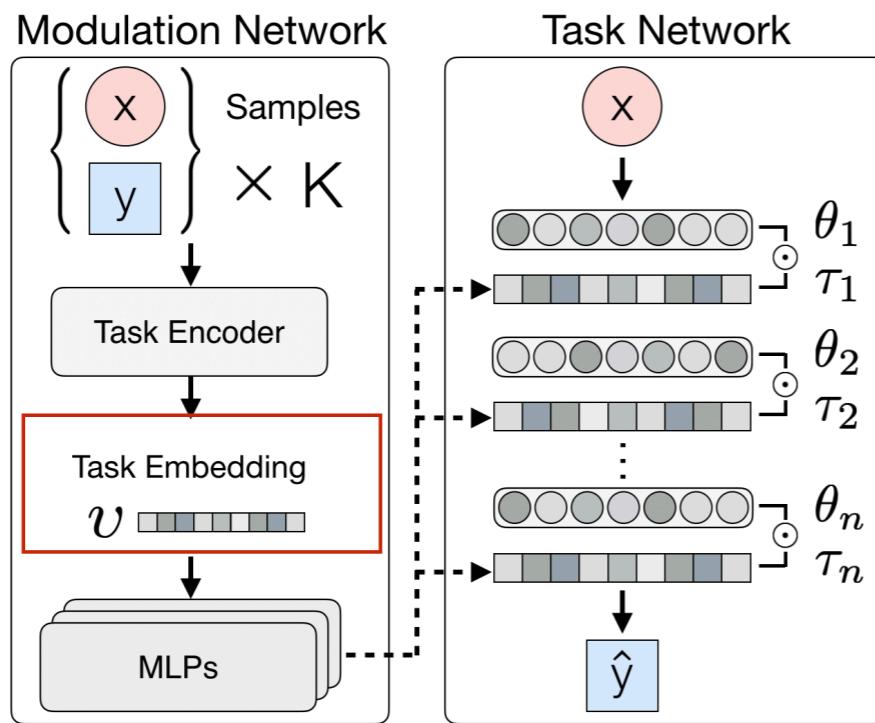
ProMP



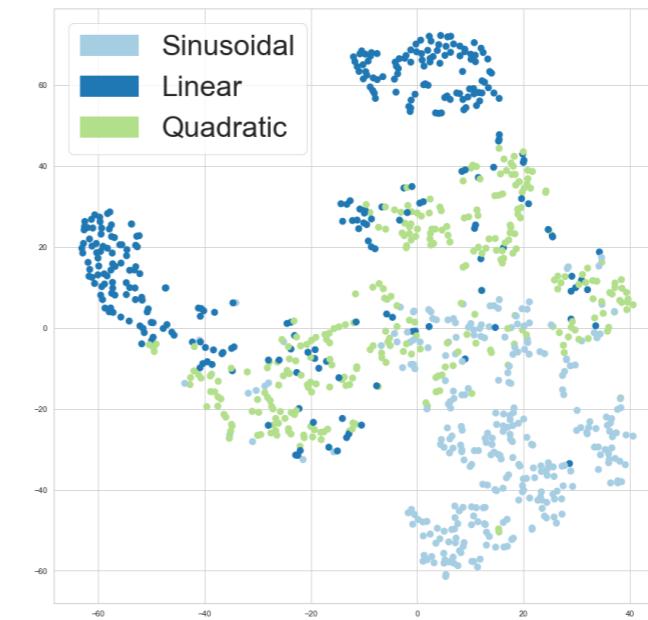
Ours



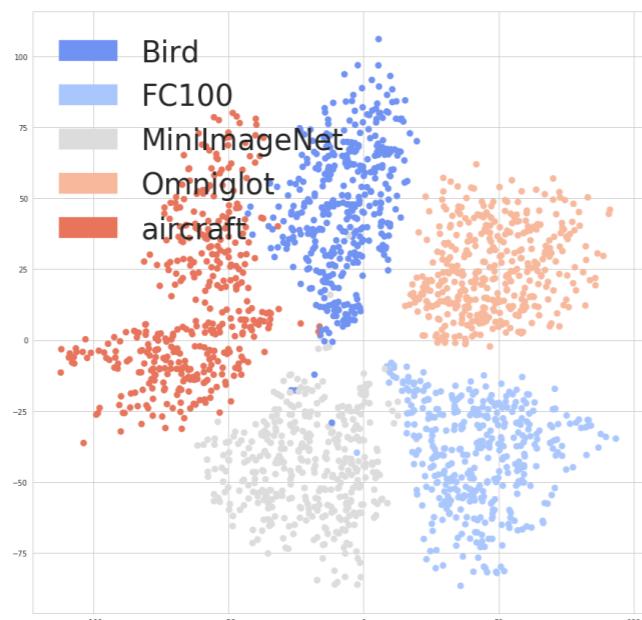
Learned Task Embedding (tSNE plot)



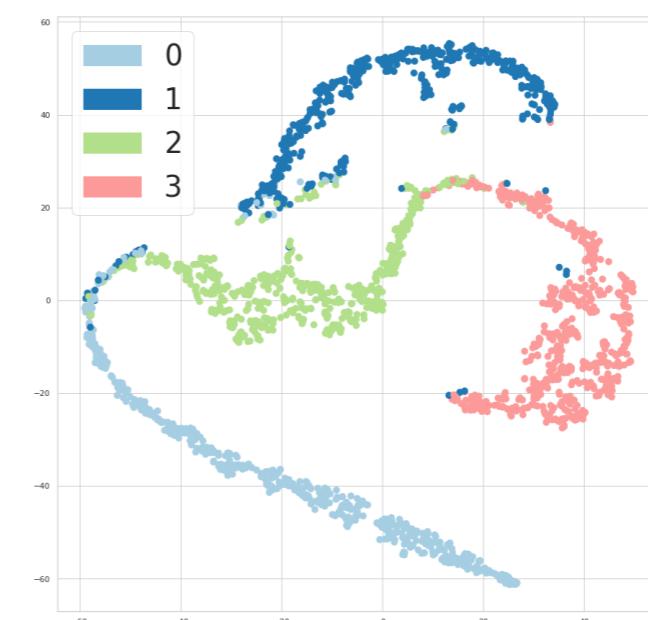
3-mode Regression



5-mode Classification

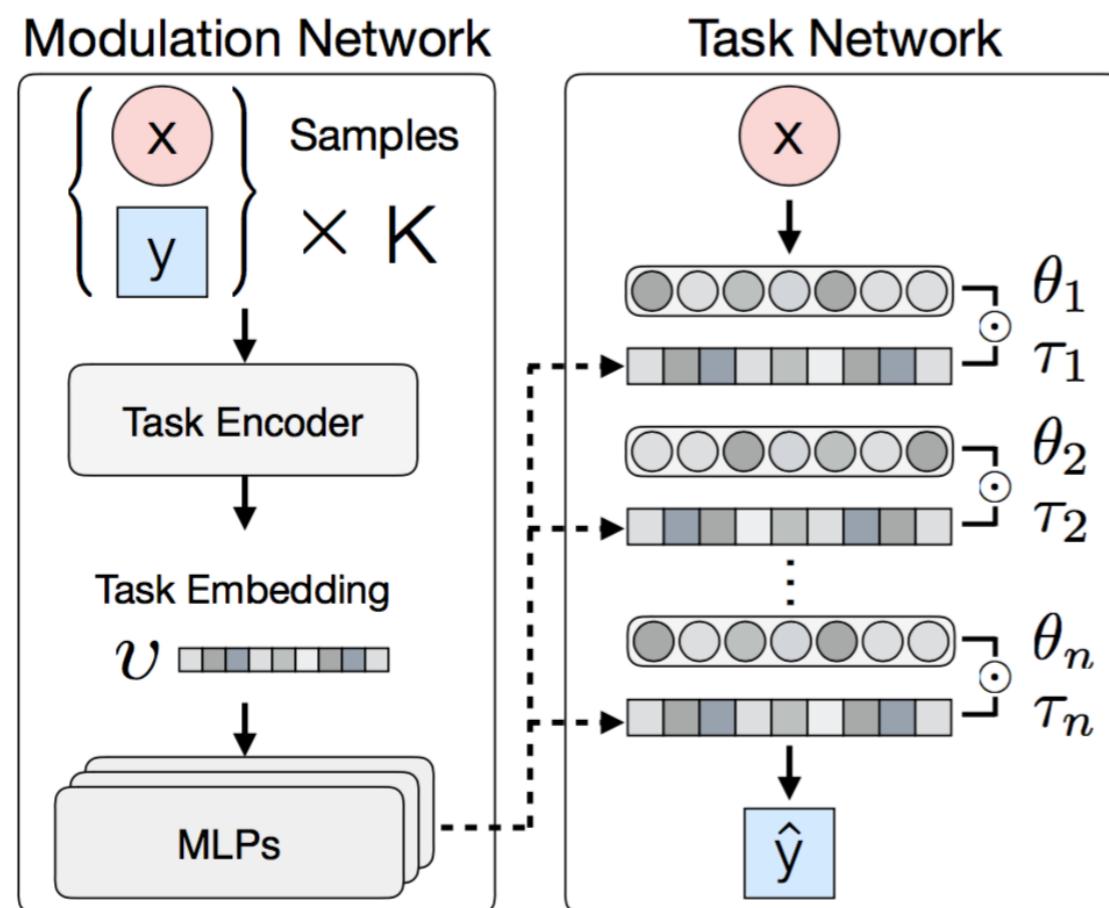


4-mode Reacher



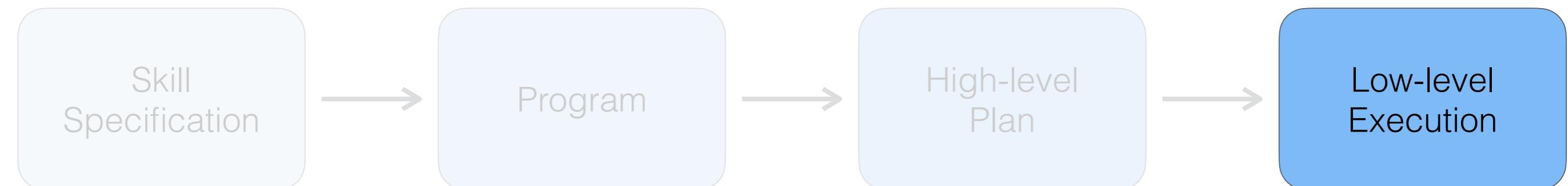
Takeaway

- MAML struggles at learning from multimodal task distributions
- We propose **multimodal MAML** to alleviate the issue

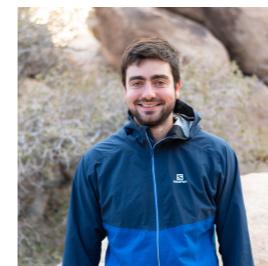


Skill-based Meta-Reinforcement Learning

Deep RL workshop @ NeurIPS 2021
Meta-learning workshop @ NeurIPS 2021
submitted to ICLR 2022



Taewook Nam



Karl Pertsch

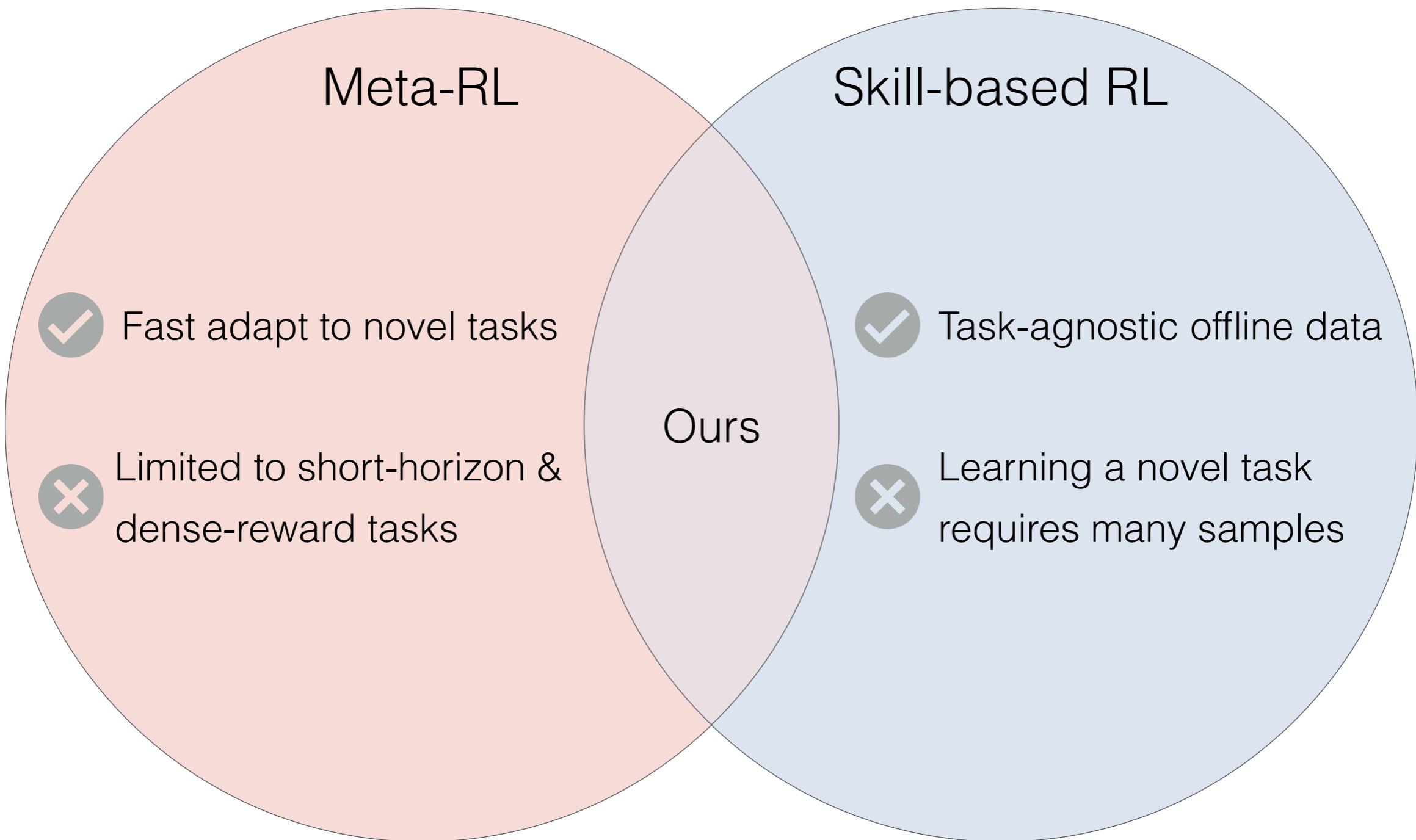


Sung Ju Hwang

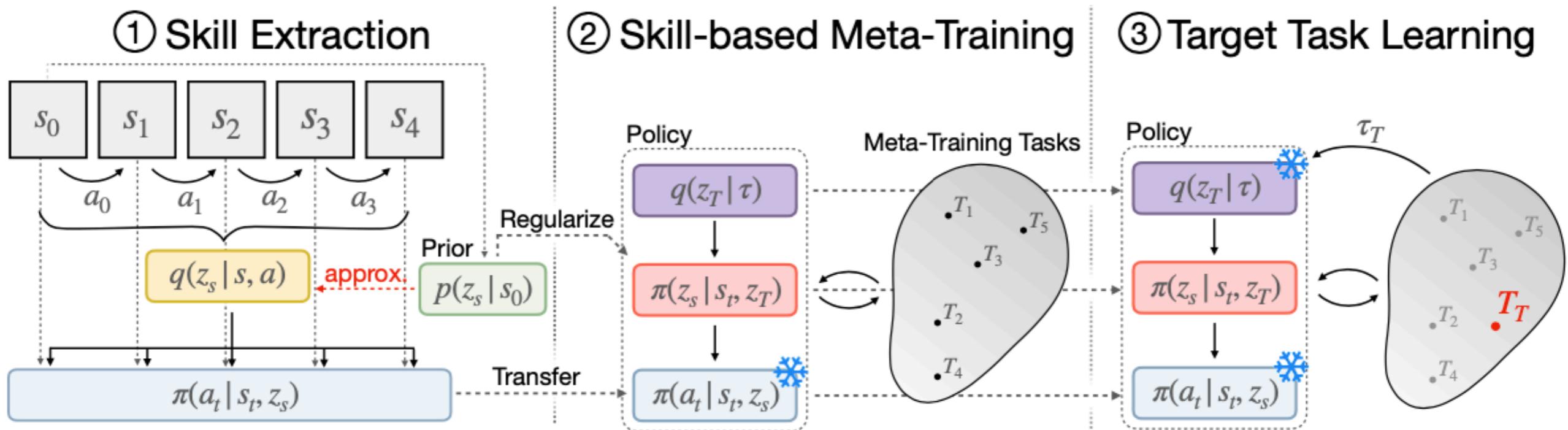


Joseph J. Lim

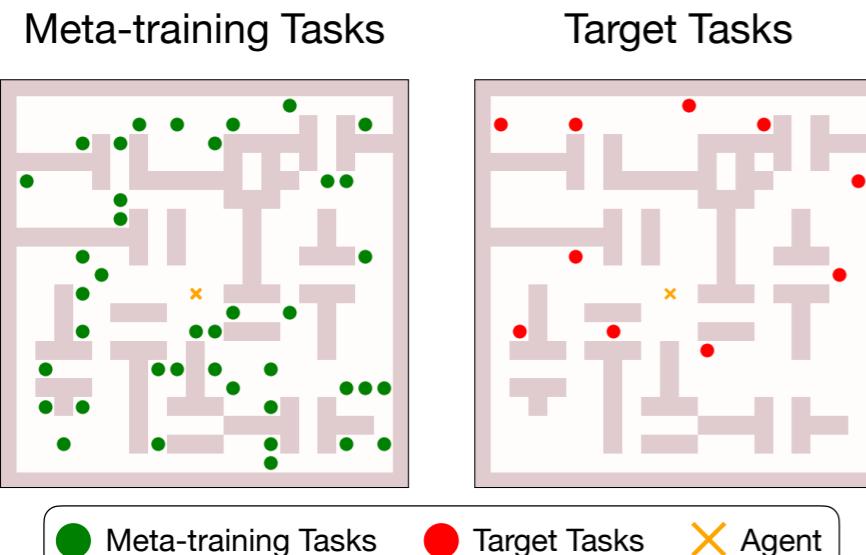
Meta-RL with Skills



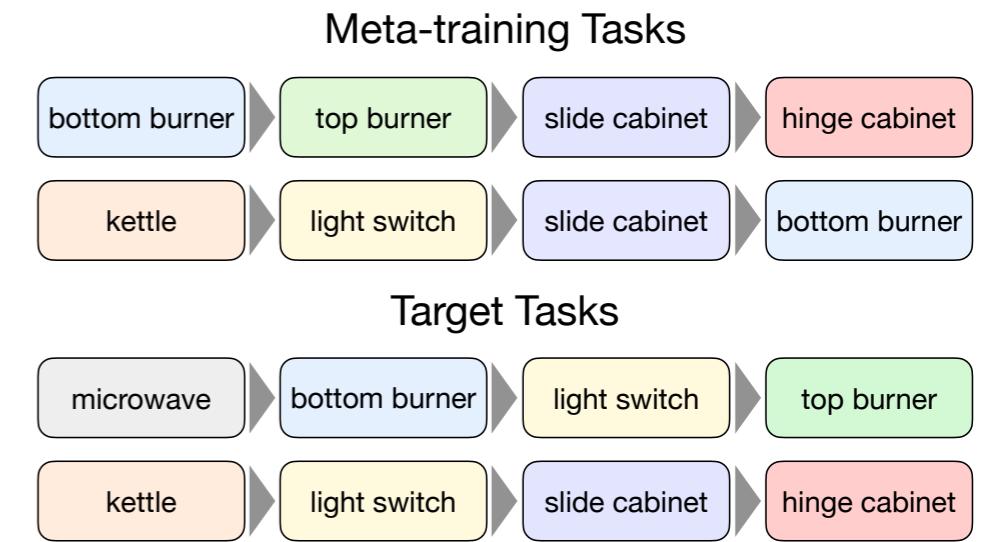
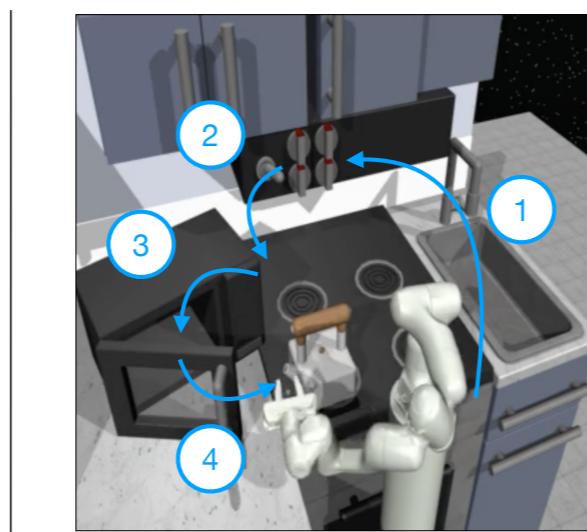
SiMPL: Skill-based Meta Policy Learning



Environments

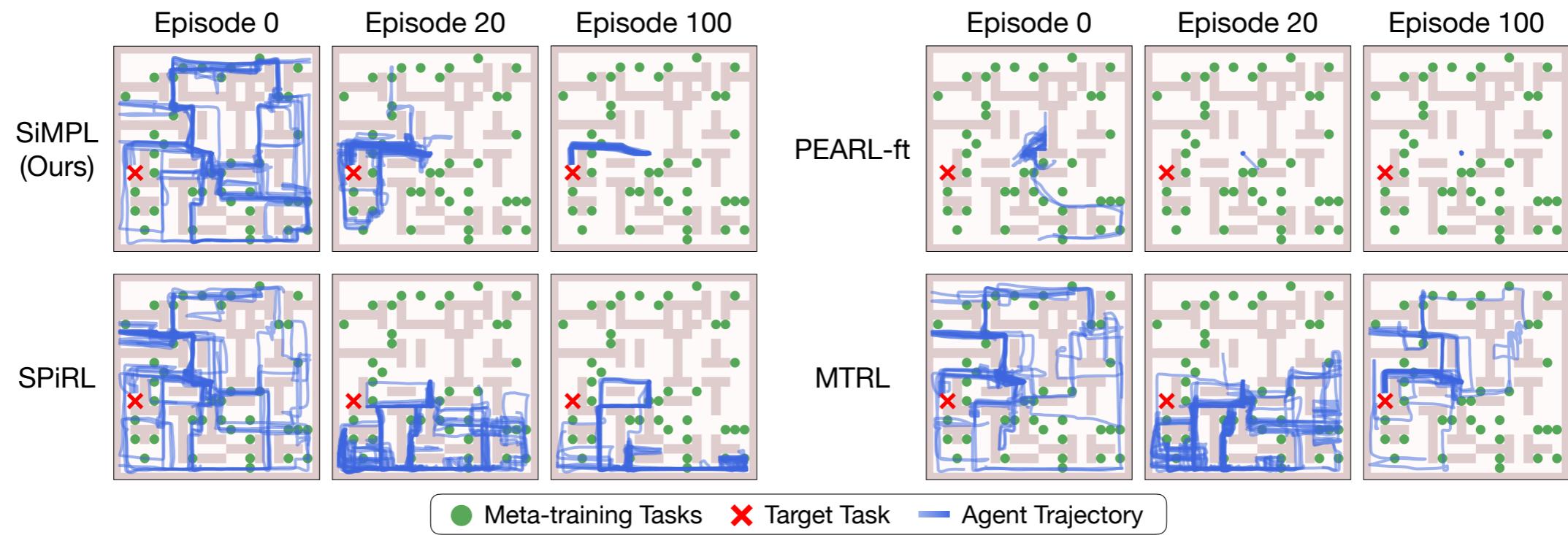
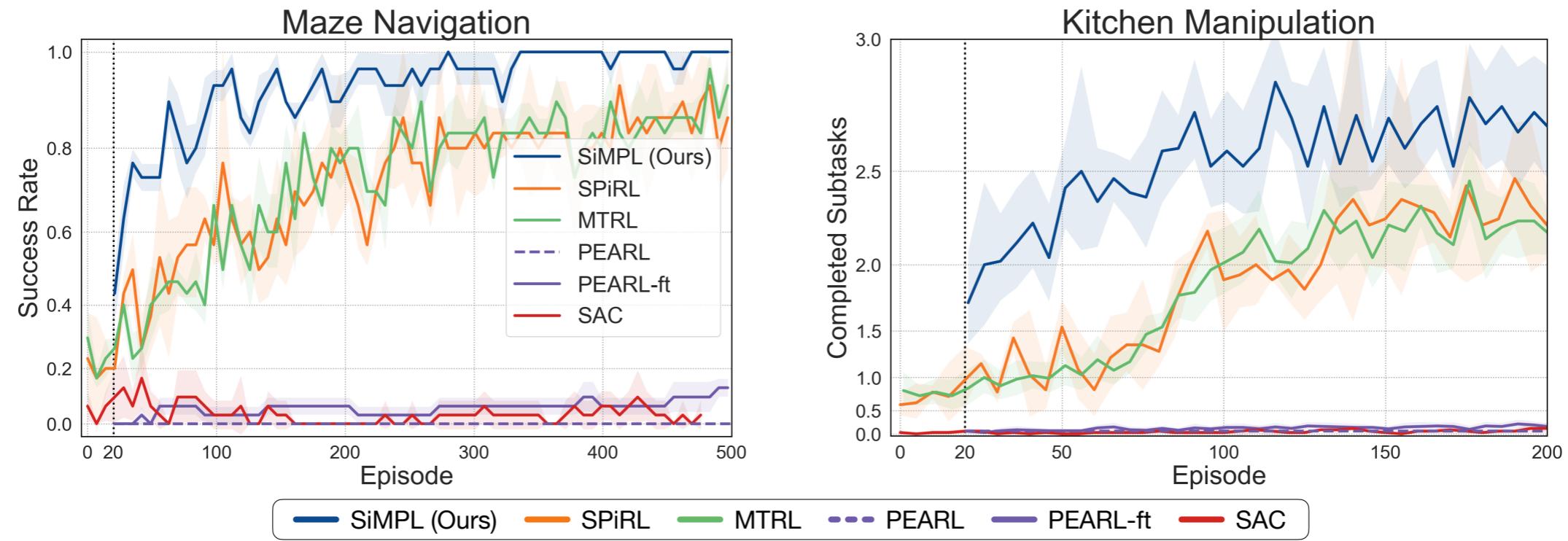


(a) Maze Navigation



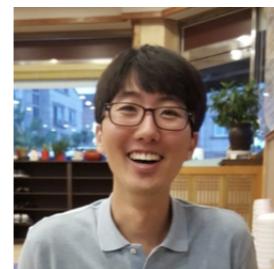
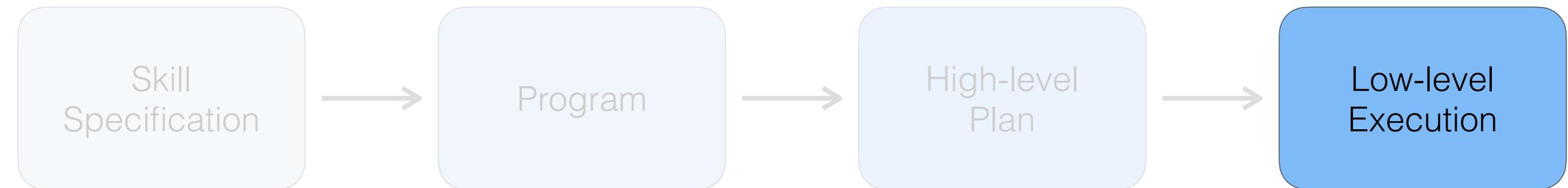
(b) Kitchen Manipulation

Results

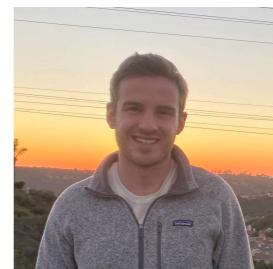


Generalizable Imitation Learning from Observation via Inferring Goal Proximity

NeurIPS 2021



Youngwoon Lee

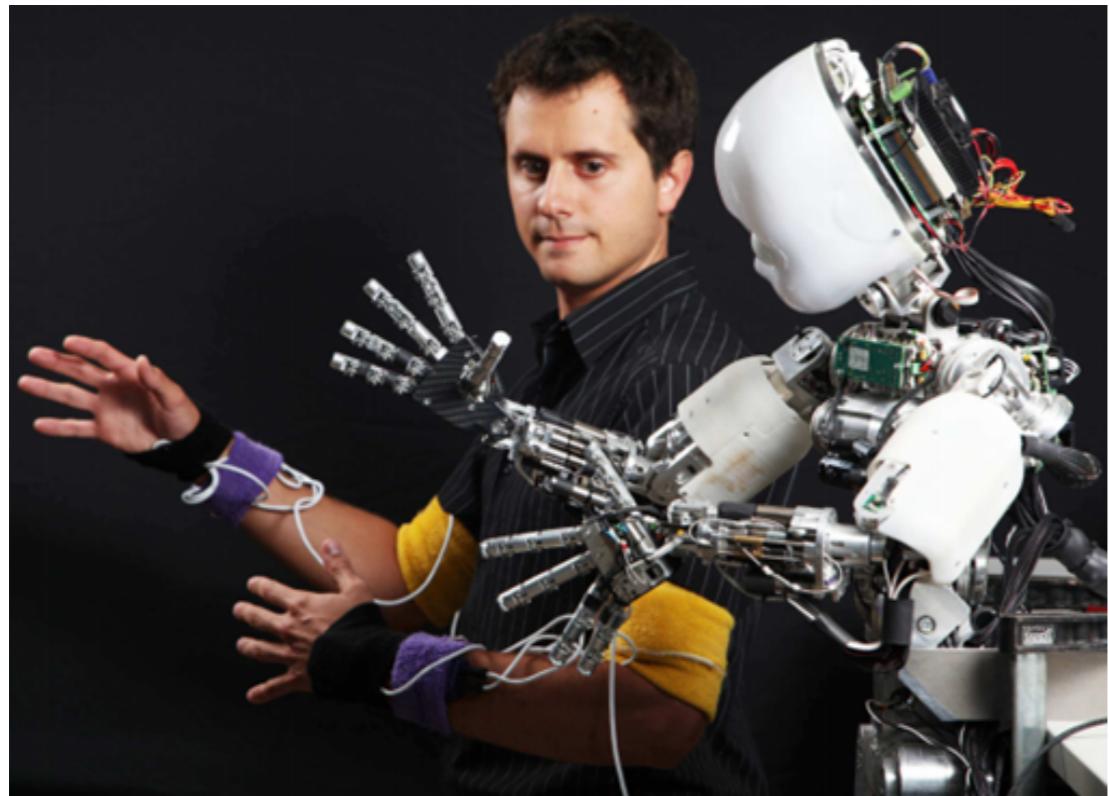


Andrew Szot



Joseph J. Lim

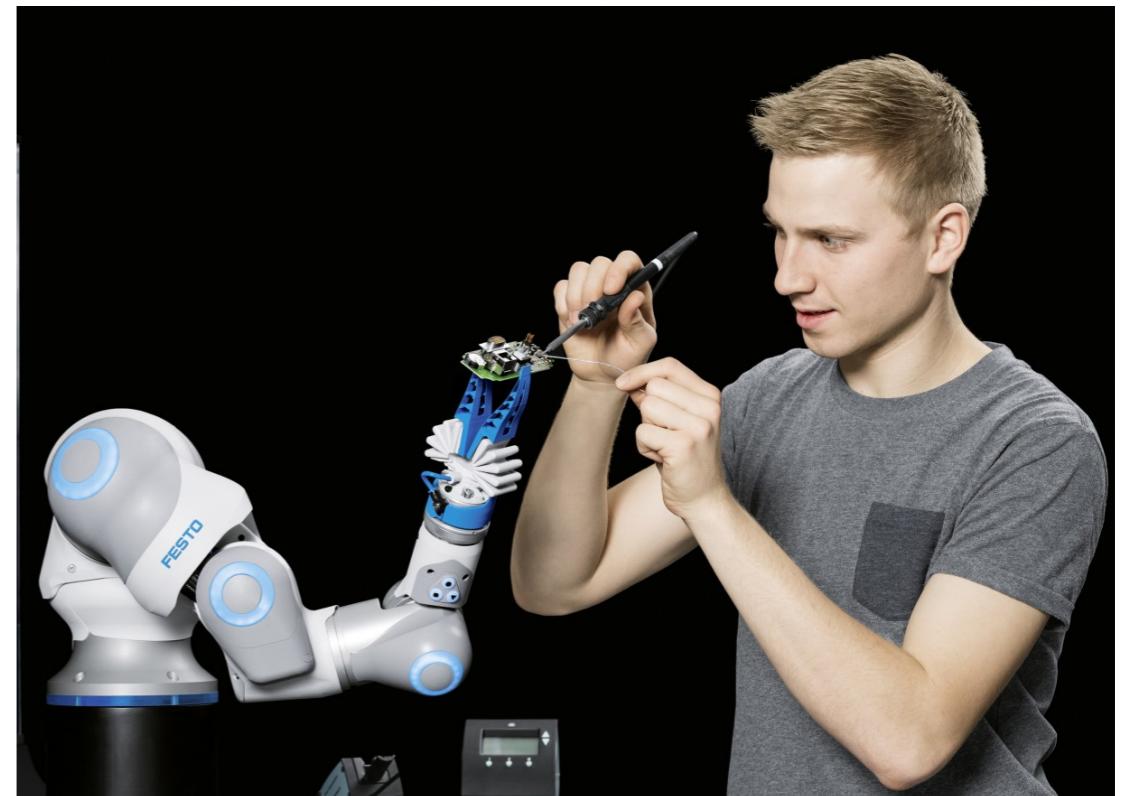
Learning from Demonstration



with expert's actions

Demo: $\{s_1, a_1, s_2, a_2, s_3, a_3, \dots\}$

Learning from Observation

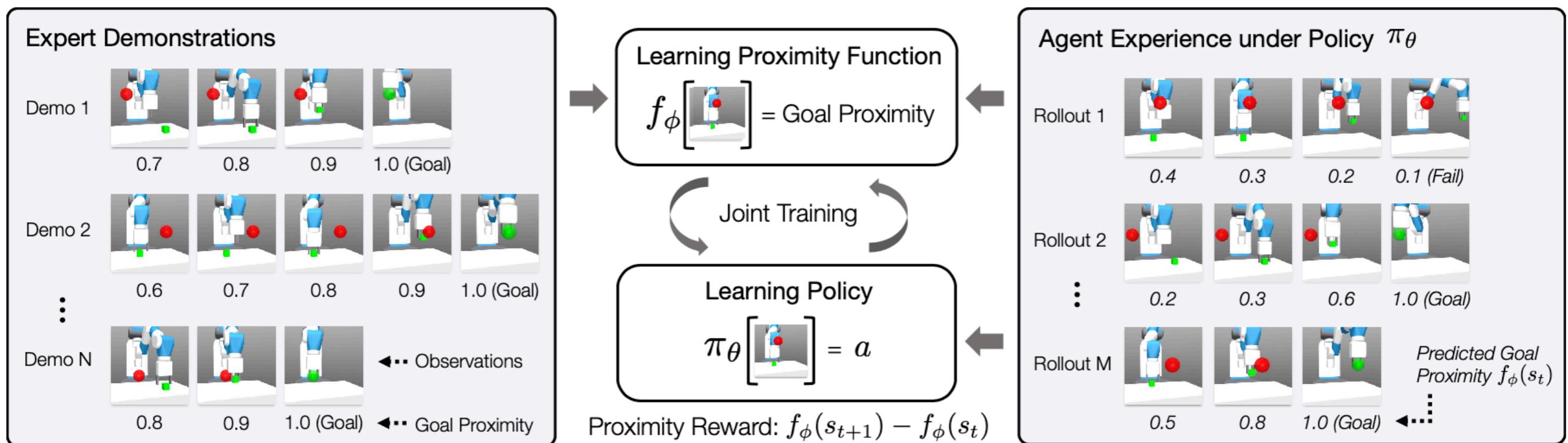


vs.

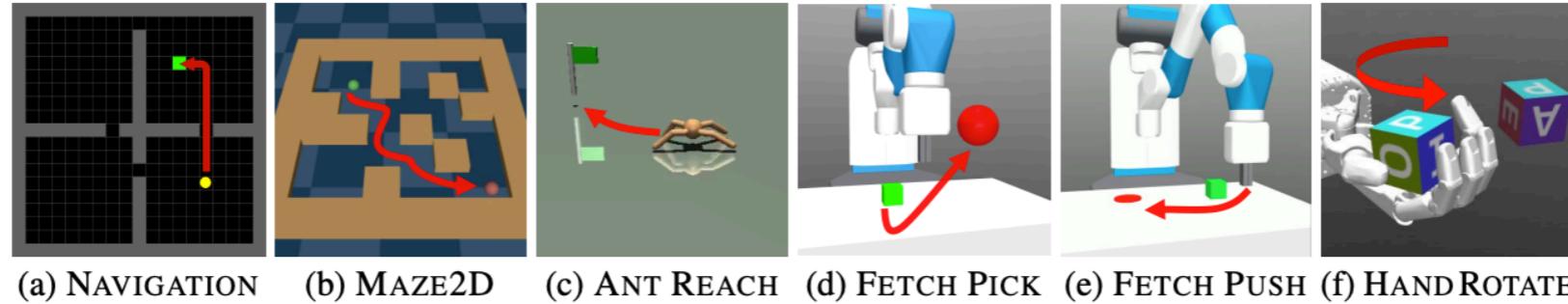
without expert's actions

Demo: $\{s_1, s_2, s_3, \dots\}$

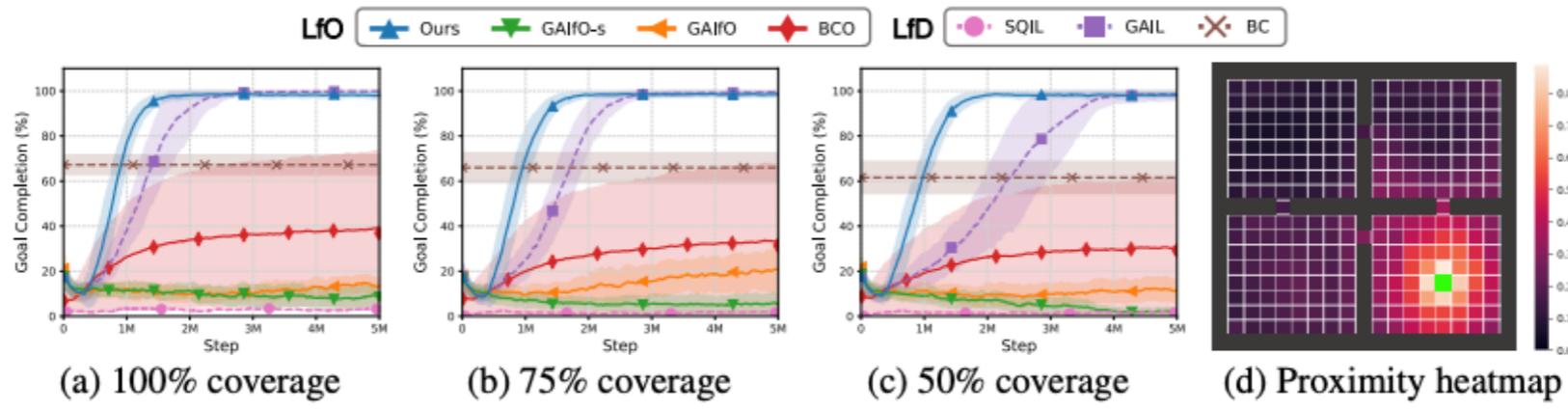
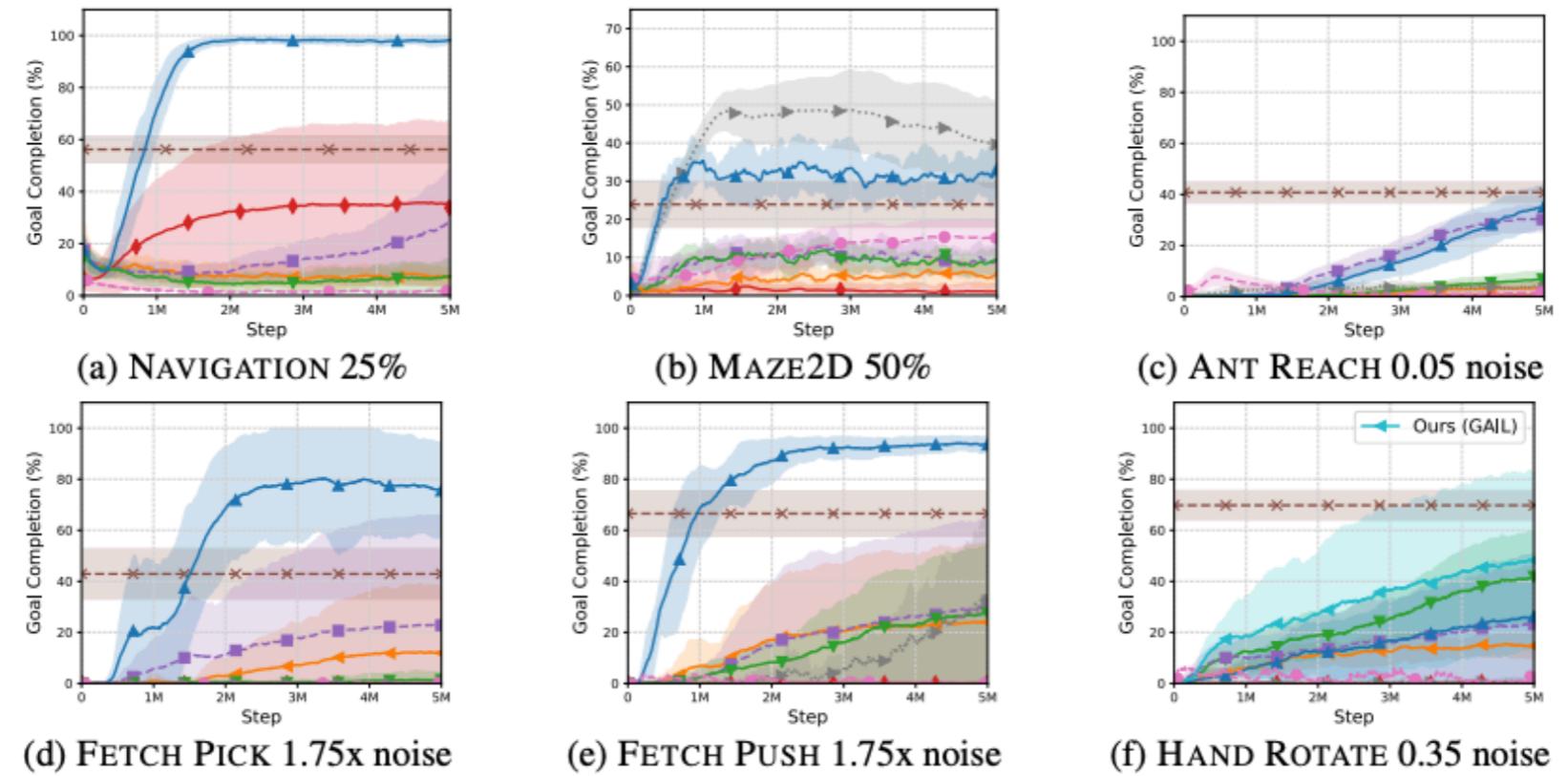
Learning from Observation via Inferring Goal Proximity



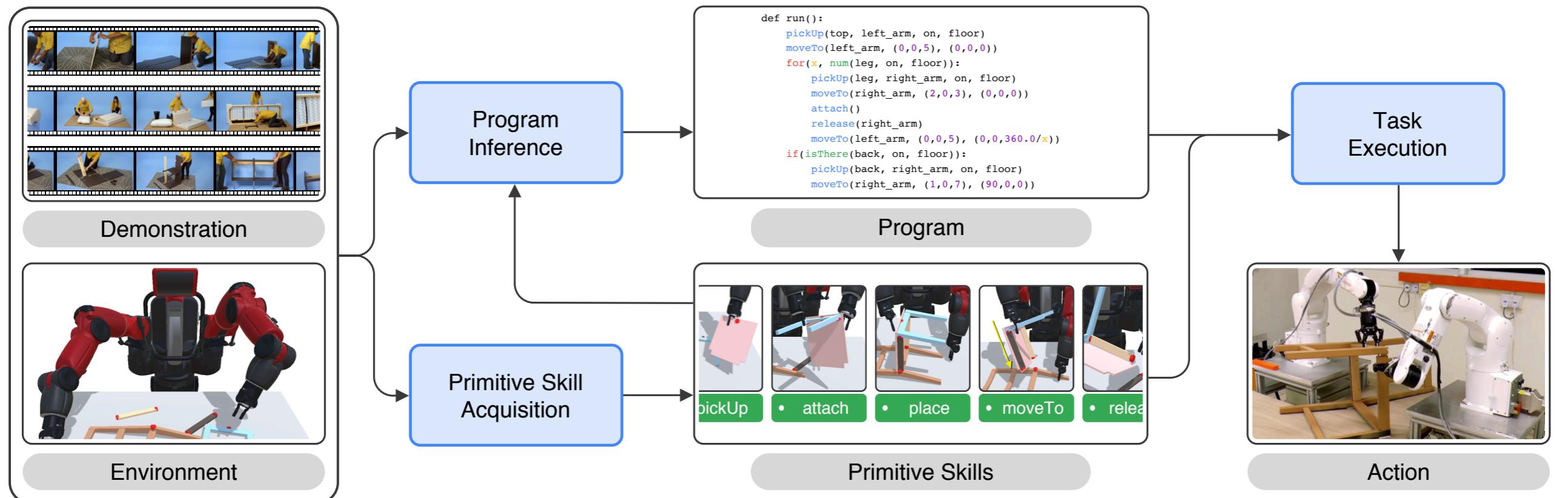
Experiments



LfO [Ours (blue star), GAIfo-s (green triangle), GAIfo (orange triangle), BCO (red diamond)] LfD [SQIL (pink circle), GAIL (purple square), BC (brown cross)] LfO+reward [GoalGAIL (grey right-pointing triangle)]



Program-Guided Framework for Interpreting and Acquiring Complex Skills



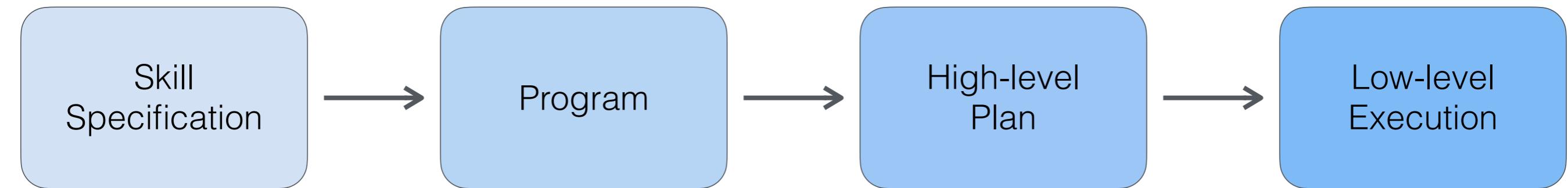
Interpretable

Programmatic / Generalizable

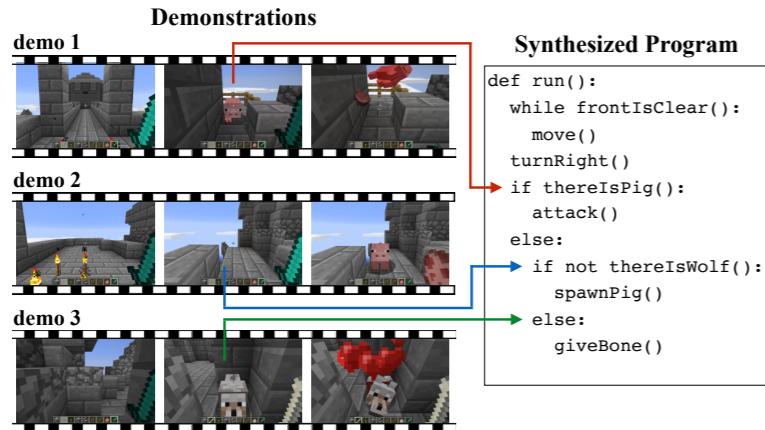
Hierarchical

Modular

Program-Guided Framework for Interpreting and Acquiring Complex Skills



Neural Program Synthesis from Diverse Demonstration Videos



Synthesized Program

```
def run():
    while frontIsClear():
        move()
        turnRight()
    if thereIsPig():
        attack()
    else:
        if not thereIsWolf():
            spawnPig()
        else:
            giveBone()
```

ICML 2018

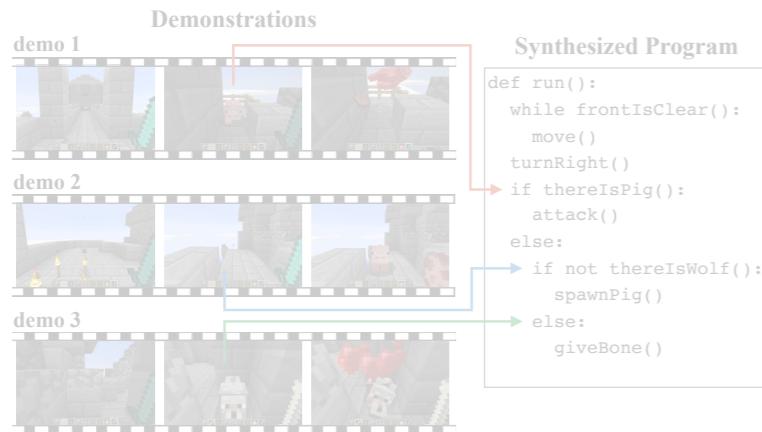
Skill
Specification

Program

High-level
Plan

Low-level
Execution

Neural Program Synthesis from Diverse Demonstration Videos



ICML 2018

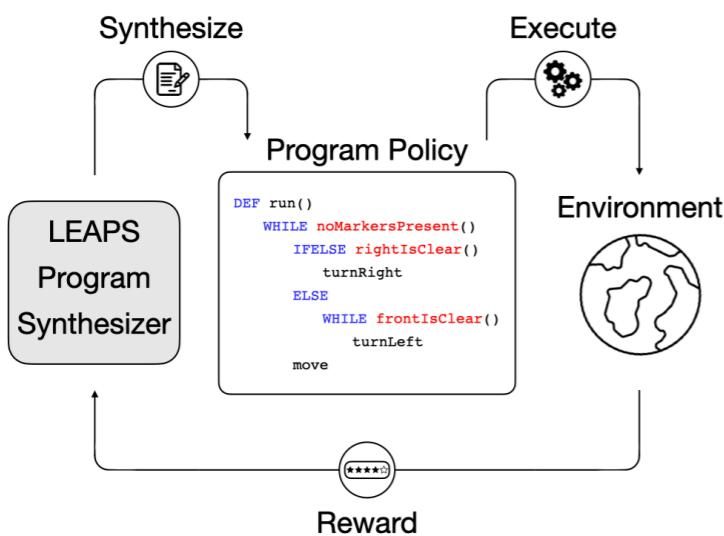
Skill Specification

Program

High-level Plan

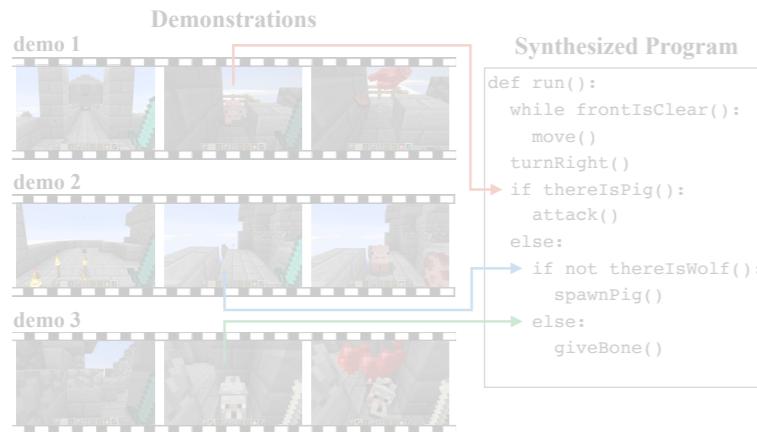
Low-level Execution

Learning to Synthesize Programs as Interpretable and Generalizable Policies



NeurIPS 2021

Neural Program Synthesis from Diverse Demonstration Videos



ICML 2018

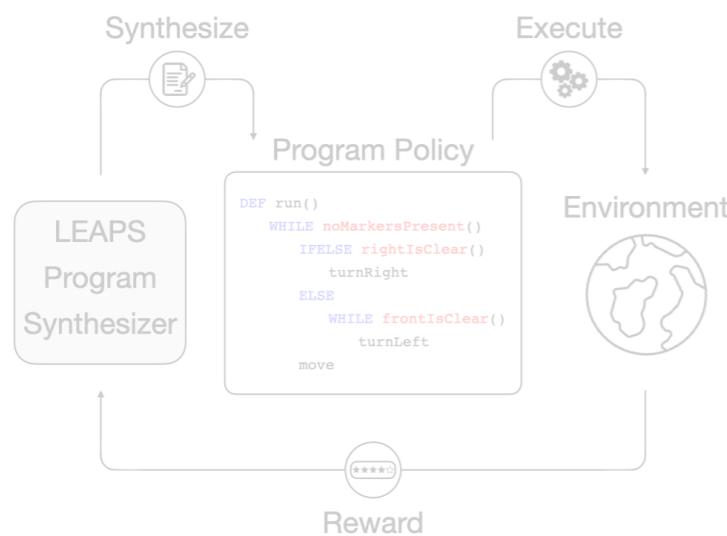
Skill Specification

Program

High-level Plan

Low-level Execution

Learning to Synthesize Programs as
Interpretable and Generalizable Policies



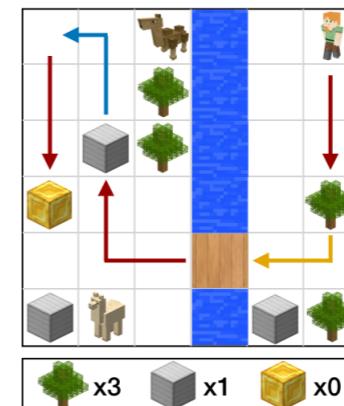
NeurIPS 2021

Program Guided Agent

Program

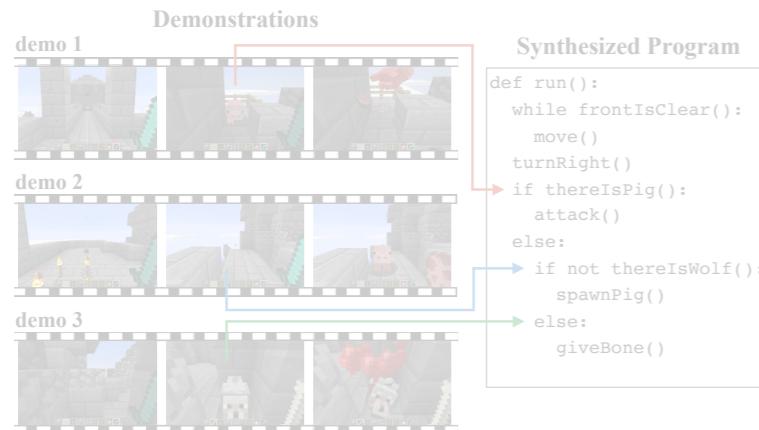
```

def Task():
    if is_there[River]:
        mine(Wood)
        build_bridge()
    if agent[Iron] < 3:
        mine(Iron)
        place(Iron, 2, 3)
    else:
        goto(4, 2)
    while env[Gold] > 0:
        mine(Gold)
    
```



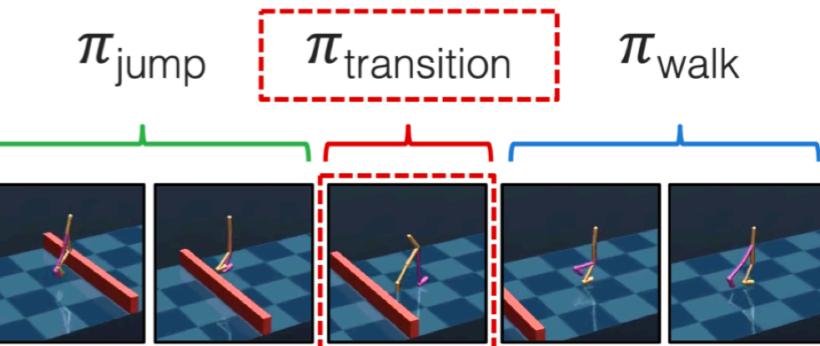
ICLR 2020 (Spotlight)

Neural Program Synthesis from Diverse Demonstration Videos



ICML 2018

Composing Complex Skills by Learning Transition Policies



ICLR 2019

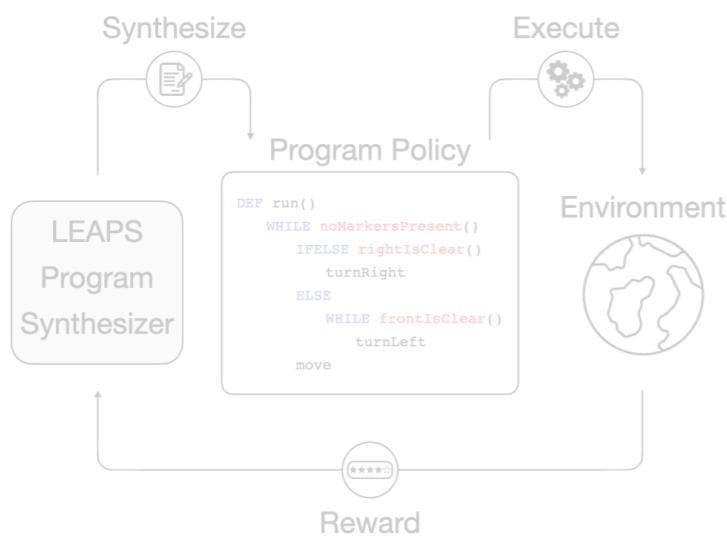
Skill Specification

Program

High-level Plan

Low-level Execution

Learning to Synthesize Programs as Interpretable and Generalizable Policies

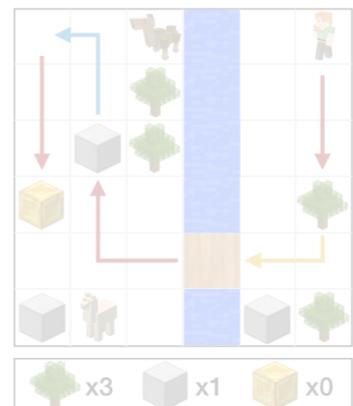


NeurIPS 2021

Program Guided Agent

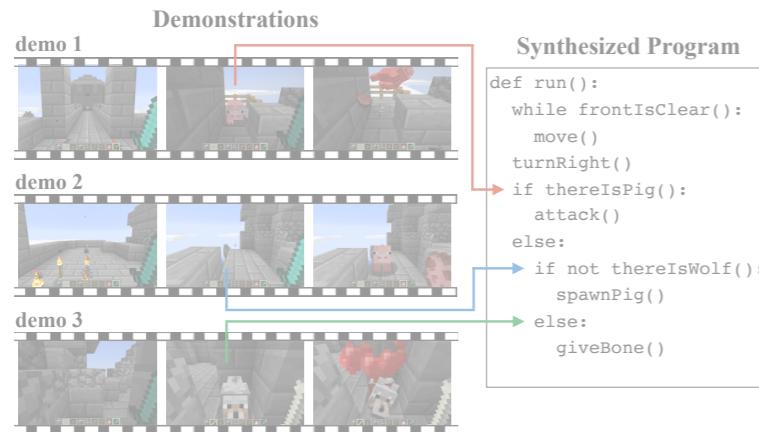
Program

```
def Task():
    if is_there[River]:
        mine(Wood)
        build_bridge()
    if agent[Iron] < 3:
        mine(Iron)
        place(Iron, 2, 3)
    else:
        goto(4, 2)
    while env[Gold] > 0:
        mine(Gold)
```

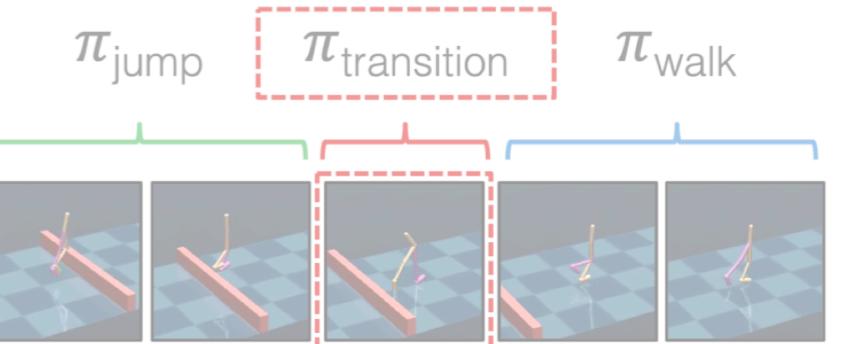


ICLR 2020 (Spotlight)

Neural Program Synthesis from Diverse Demonstration Videos



ICML 2018



ICLR 2019

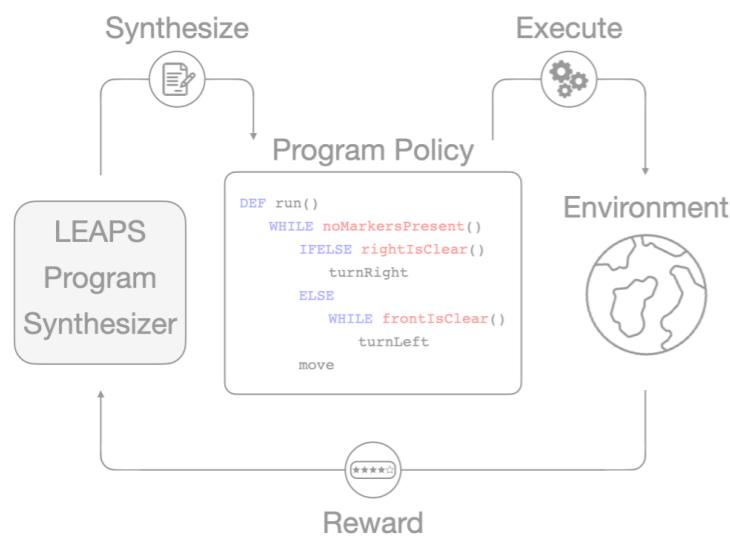
Skill
Specification

Program

High-level
Plan

Low-level
Execution

Learning to Synthesize Programs as Interpretable and Generalizable Policies

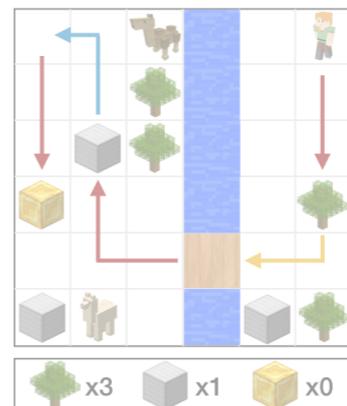


NeurIPS 2021

Program Guided Agent

Program

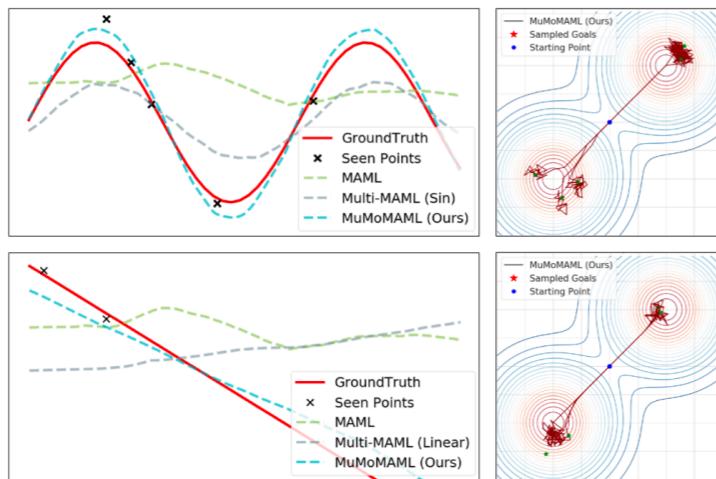
```
def Task():
    if is_there[River]:
        mine(Wood)
        build_bridge()
    if agent[Iron] < 3:
        mine(Iron)
        place(Iron, 2, 3)
    else:
        goto(4, 2)
    while env[Gold] > 0:
        mine(Gold)
```



ICLR 2020 (Spotlight)

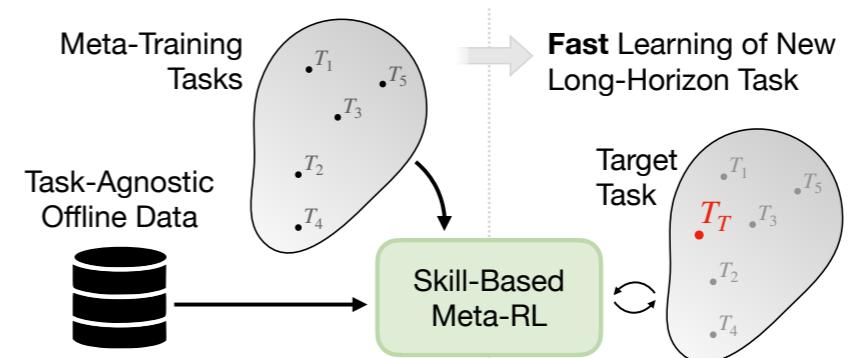
Low-level Execution

Toward Multimodal Model-Agnostic Meta-Learning



Meta-learning workshop @ NeurIPS 2018

Skill-based Meta-Reinforcement Learning

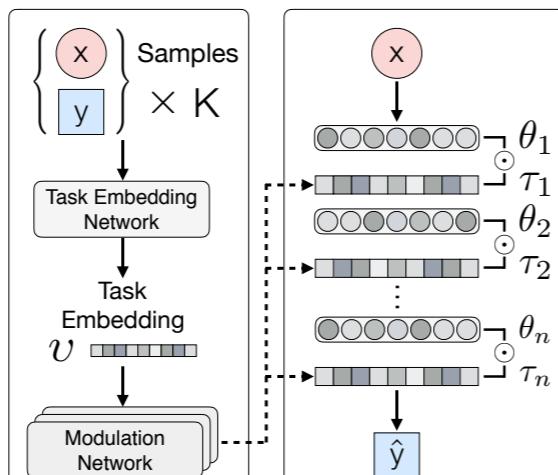


Deep RL workshop @ NeurIPS 2021

Meta-learning workshop @ NeurIPS 2021

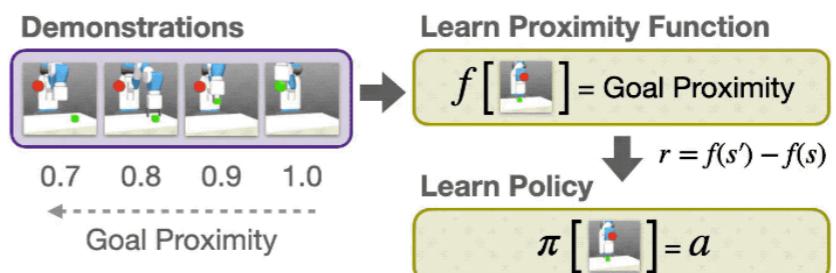
submitted to ICLR 2022

Multimodal Model-Agnostic Meta-Learning via Task-Aware Modulation



NeurIPS 2019 (Spotlight)

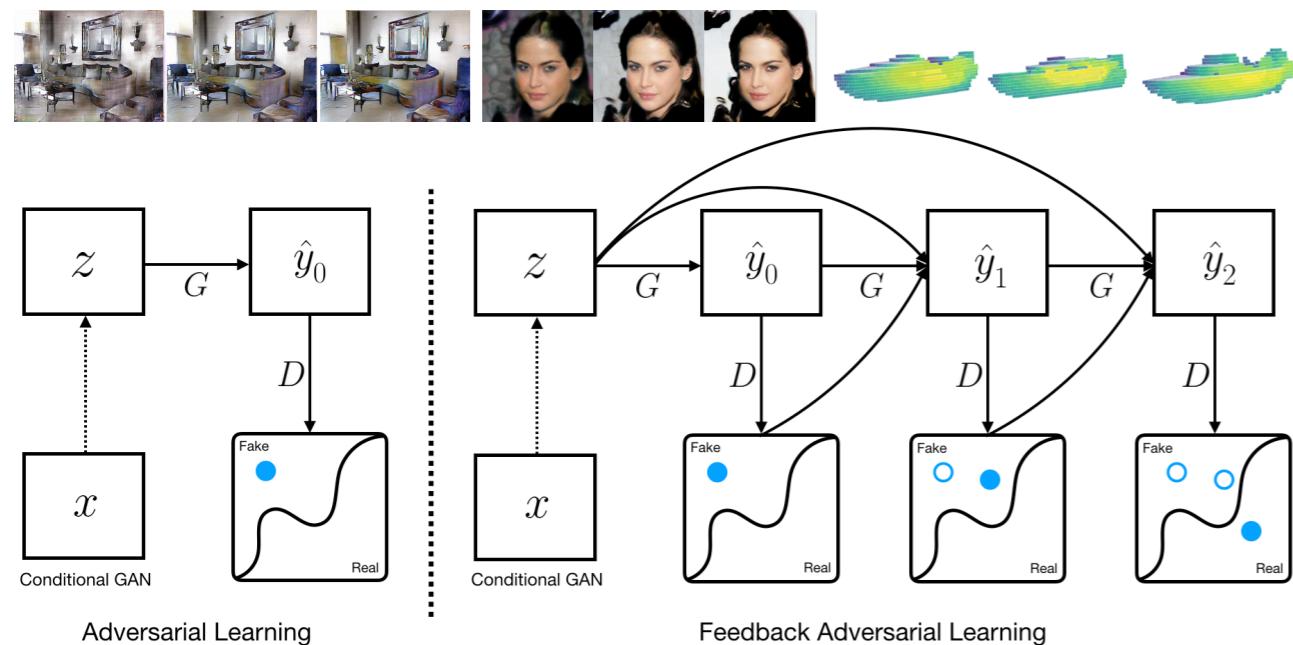
Generalizable Imitation Learning from Observation via Inferring Goal Proximity



NeurIPS 2021

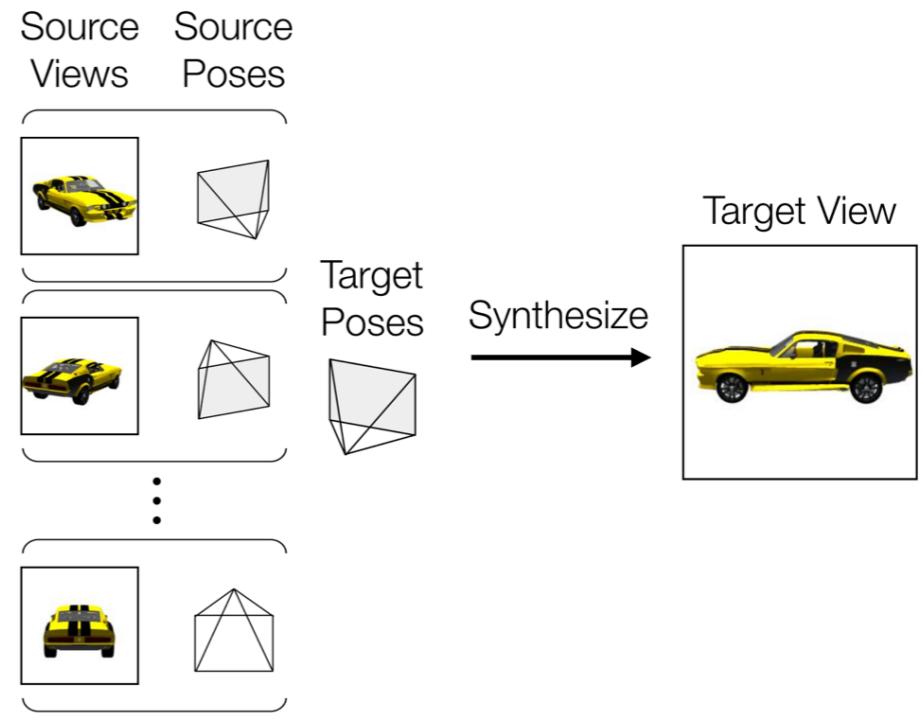
Computer Vision / 3D Vision / Image Synthesis / GAN / Medical Imaging

Feedback Adversarial Learning: Spatial Feedback for Improving Generative Adversarial Networks



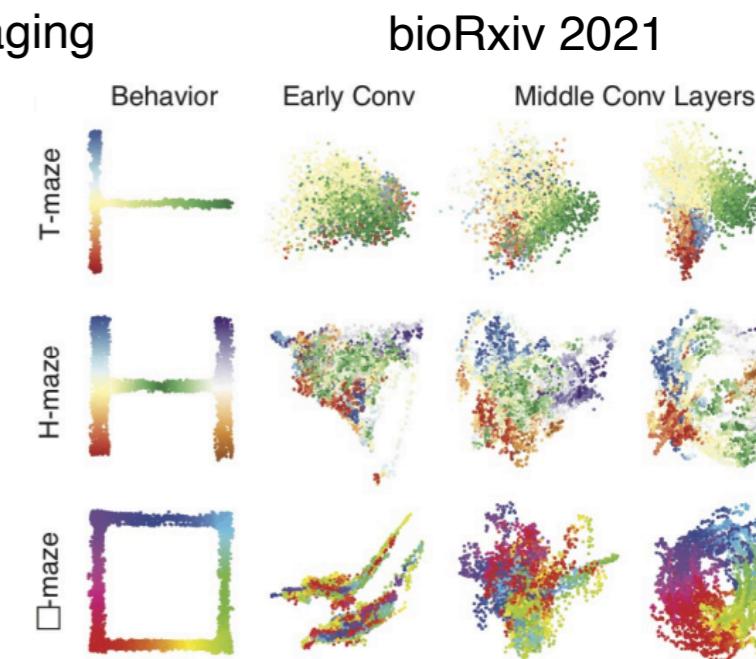
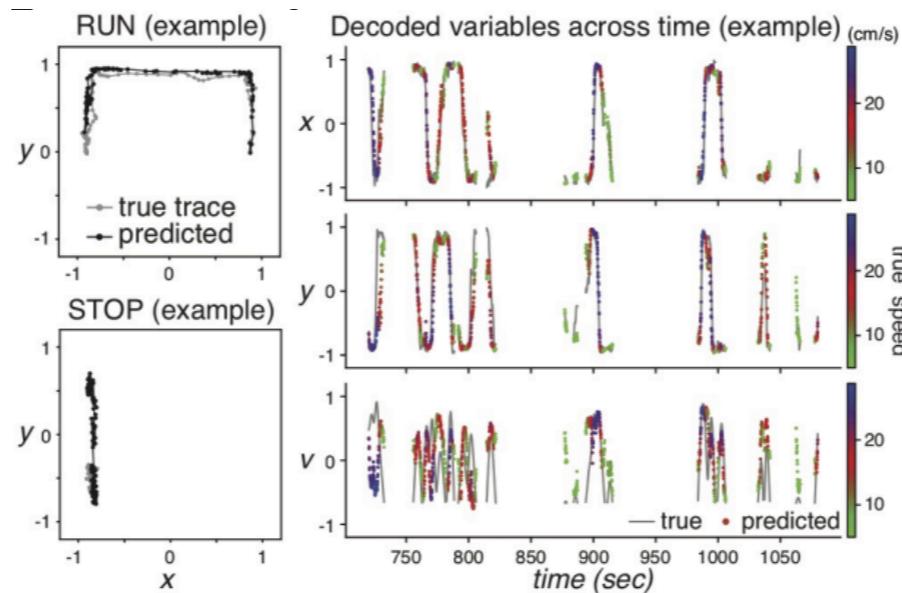
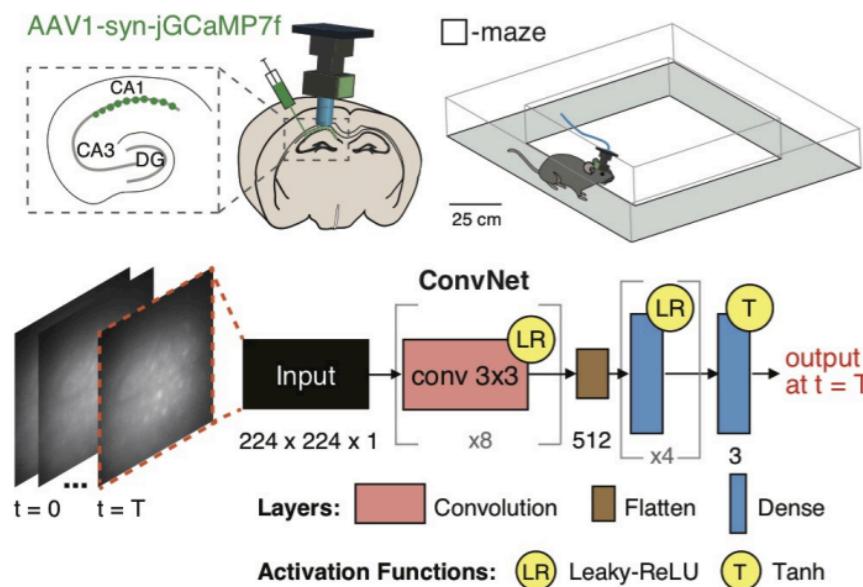
CVPR 2019

Multi-view to Novel view: Synthesizing Views with Self-Learned Confidence



ECCV 2018

Behavioral Clusters Revealed by End-to-end Decoding from Microendoscopic Imaging



bioRxiv 2021



Thank You

Questions?