# Neural Programmer-Interpreters

Scott Reed and Nando de Freitas
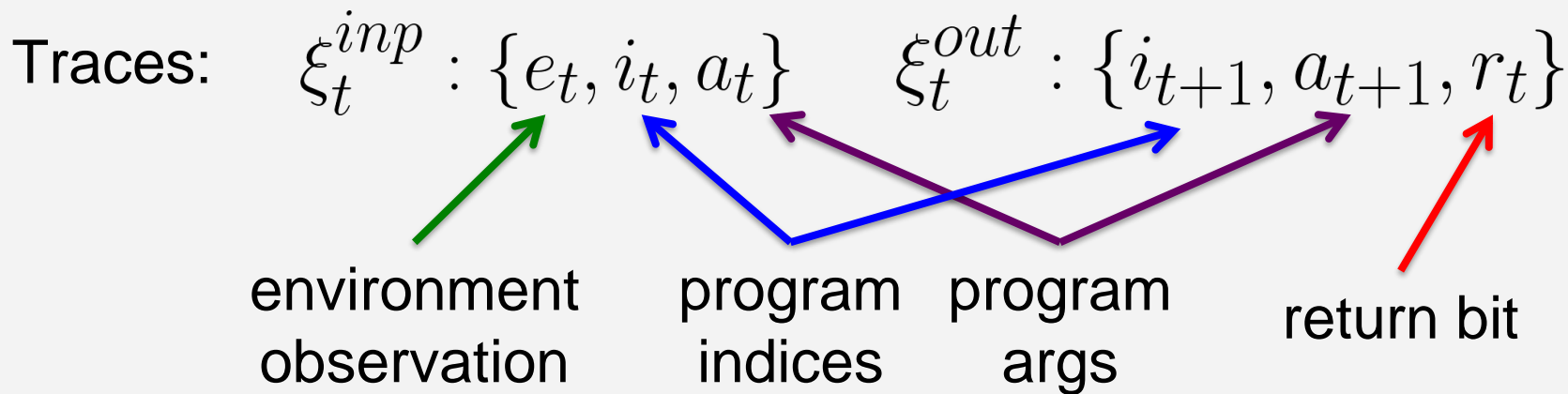
Google DeepMind

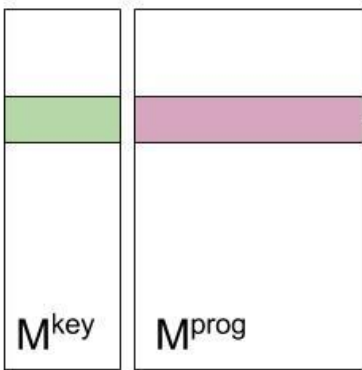# Neural Programmer Interpreter (NPI) goals:

1. **Long-term prediction:** Model potentially long sequences of actions by exploiting compositional structure.

2. **Continual learning:** Learn new programs by composing previously-learned programs, rather than from scratch.

3. **Data efficiency:** Learn generalizable programs from a small number of example traces.

4. **Interpretability**: By looking at NPI's generated commands, we can understand what it is doing at multiple levels of temporal abstraction.
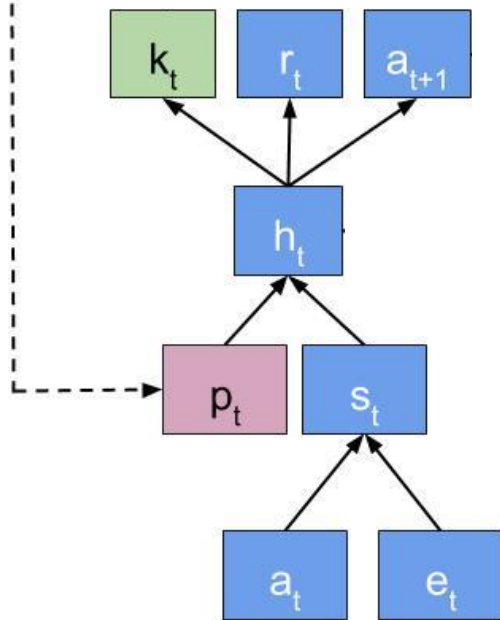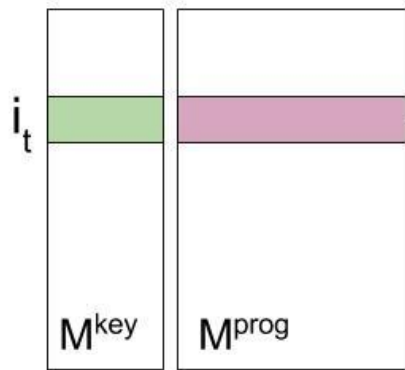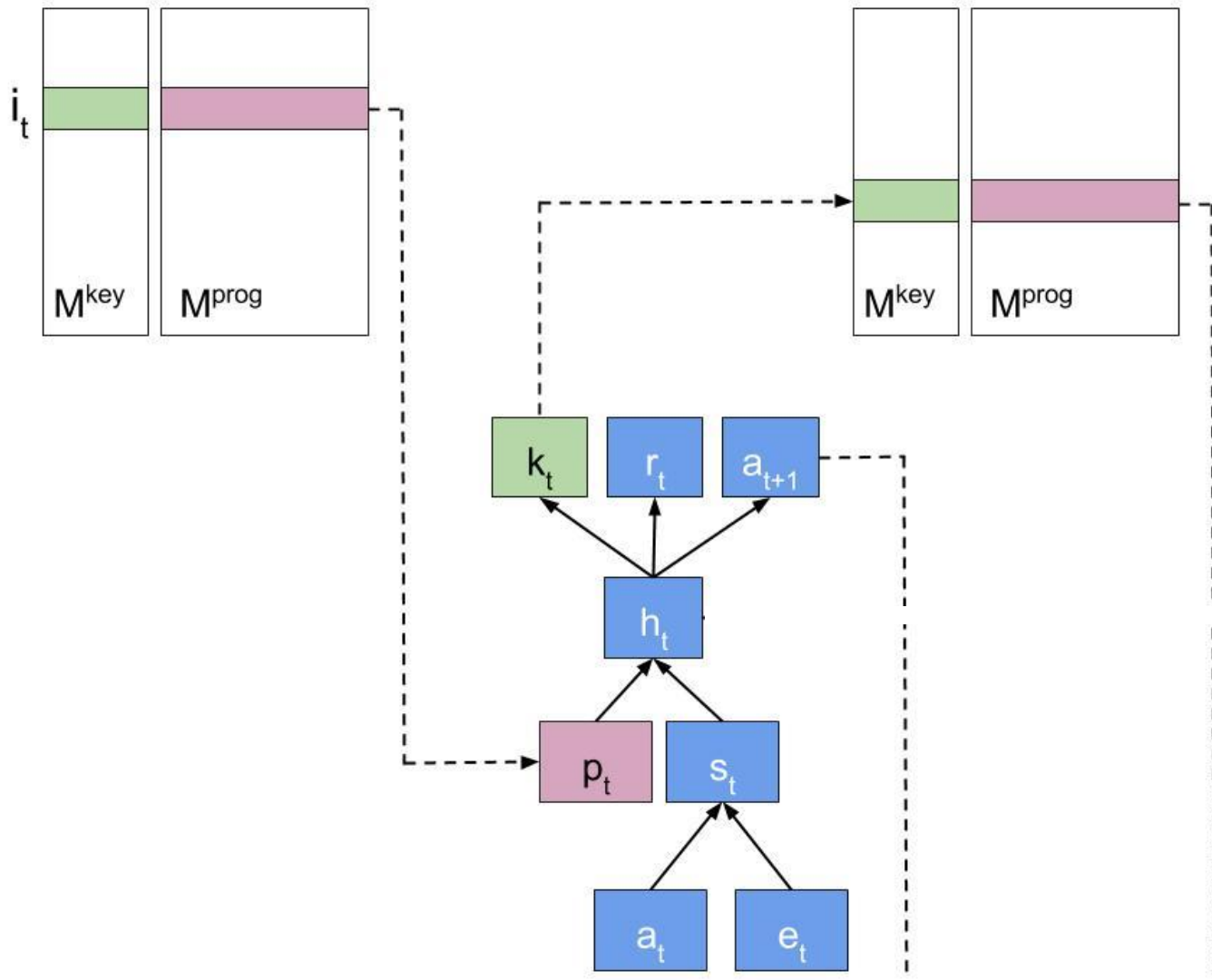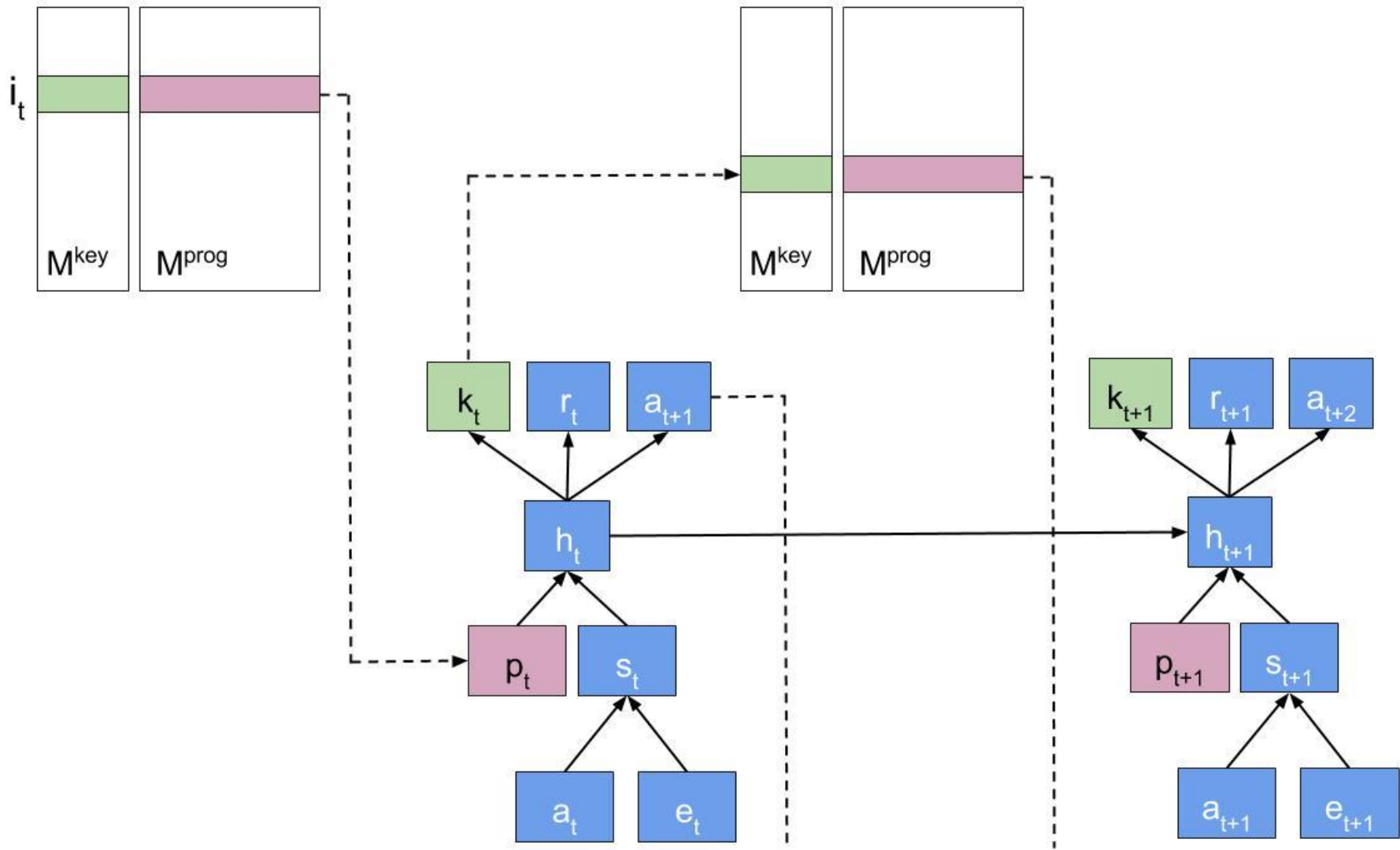
# Model

# NPI training data

Traces:

$$\xi_t^{inp} : \{e_t, i_t, a_t\} \qquad \xi_t^{out} : \{i_{t+1}, a_{t+1}, r_t\}$$



environment
observation

program
indices

program
args

return bit

$i_t$

$M^{key}$  $M^{prog}$

# Demos

# Adding numbers together - environment

|        | 0 | 0 | 0 | 9 | 6 |
|--------|---|---|---|---|---|
| input 1 | | | | | |
| input 2 | 0 | 0 | 1 | 2 | 5 |
| carry | 0 | 0 | 1 | 1 | 1 |
| output | 0 | 0 | 0 | 2 | 1 |

Addition environment interface:
- Scratch pad with the two numbers to be added, a carry row and output row.
- LEFT, RIGHT programs that can move a pointer left or right, respectively.
- WRITE program that writes a specified value to the location of a specified pointer.
- 4 read/write pointers; one per row.

# Adding numbers together – learned programs

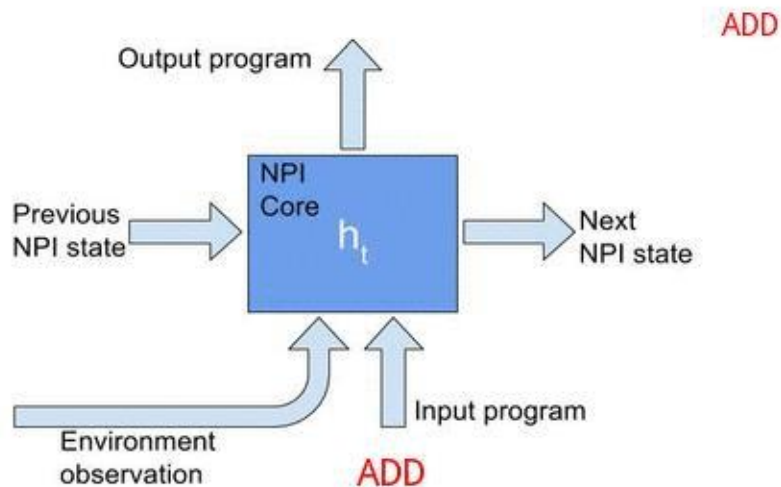| Program | Description | Calls |
|---------|-------------|-------|
| ADD | Multi-digit addition | ADD1, LSHIFT |
| ADD1 | Single-digit add | CARRY, ACT |
| CARRY | Mark a 1 in the carry row 1 step left. | ACT, LSHIFT, RSHIFT |
| LSHIFT | Shift specified pointer 1 step left. | ACT |
| RSHIFT | Shift specified pointer 1 step right. | ACT |
| ACT | Move pointer or write to the scratch pad | - |

# Adding numbers together

# Bubble sort - environment

array
t=0

| 3 | 2 | 4 | 9 | 1 |

t=1

| 3 | 2 | 4 | 9 | 1 |

t=2

| 2 | 3 | 4 | 9 | 1 |

t=3

| 2 | 3 | 4 | 9 | 1 |

Sorting environment interface:
- Scratch pad with the array to be sorted.
- Read/write pointers.
- LEFT, RIGHT programs that can move a specified pointer left or right, respectively.
- SWAP program that swaps the values at two specified pointer locations.

# Bubble sort – learned programs

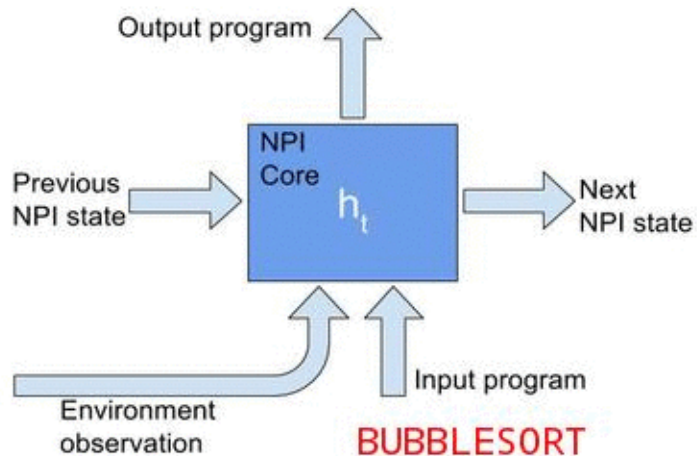| Program | Description | Calls |
|---------|-------------|-------|
| BUBBLESORT | Sort numbers in ascending order | BUBBLE, RESET |
| BUBBLE | Perform one sweep of bubble sort | BSTEP, ACT |
| RESET | Move pointers all back to the left | LSHIFT |
| BSTEP | Conditionally swap and advance pointers | COMPSWAP, RSHIFT |
| COMPSWAP | Conditionally swap two pointer values | ACT |
| LSHIFT | Shift specified pointer 1 step left. | ACT |
| RSHIFT | Shift specified pointer 1 step right. | ACT |
| ACT | Perform a swap or move a pointer. | - |

# Bubble sort



**Input array**

**NPI inference**

**Generated commands**

BUBBLESORT

*
*
*

Output program

NPI Core $h_t$

Previous NPI state

Next NPI state

Environment observation

Input program

BUBBLESORT

# 3D car models - environment



3D cars environment interface:
- Rendering of the car (pixels).
- Target angle and elevation coordinates.
- LEFT, RIGHT, UP, DOWN programs that can move the car 15 degrees at a time.
- **The current car pose is NOT provided**.
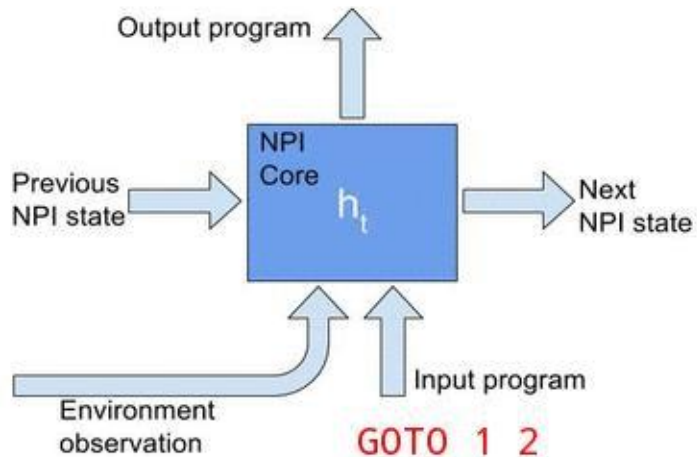
# 3D car models – learned programs

| Program | Description | Calls |
|---------|-------------|-------|
| GOTO | Change 3D car pose to match target | HGOTO, VGOTO |
| HGOTO | Move horizontally to target angle | LGOTO, RGOTO |
| LGOTO | Move left to target | ACT |
| RGOTO | Move right to target | ACT |
| VGOTO | Move vertically to target elevation | UGOTO, DGOTO |
| UGOTO | Move up to target | ACT |
| DGOTO | Move down to target | ACT |
| ACT | Move 15 degrees up, down, left or right. | - |

# Canonicalizing the view of 3D car models
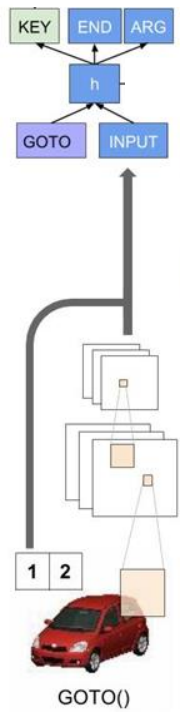


**Car rendering**

**NPI inference**

Output program

Previous NPI state → NPI Core $h_t$ → Next NPI state
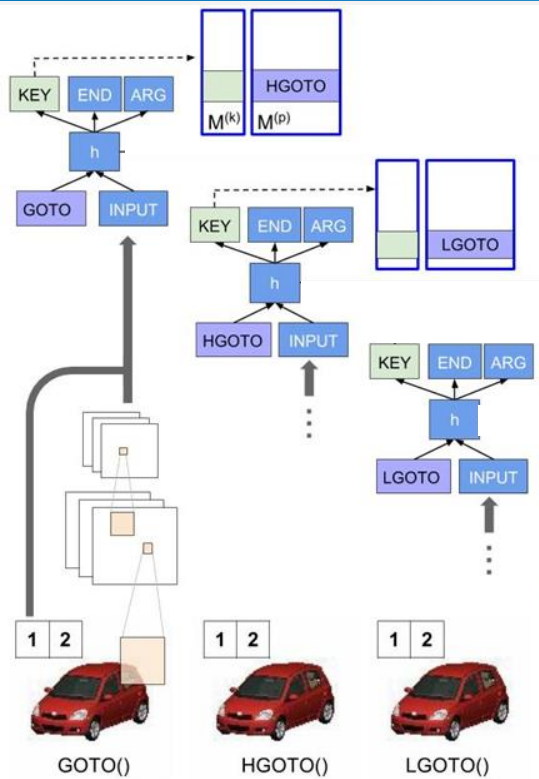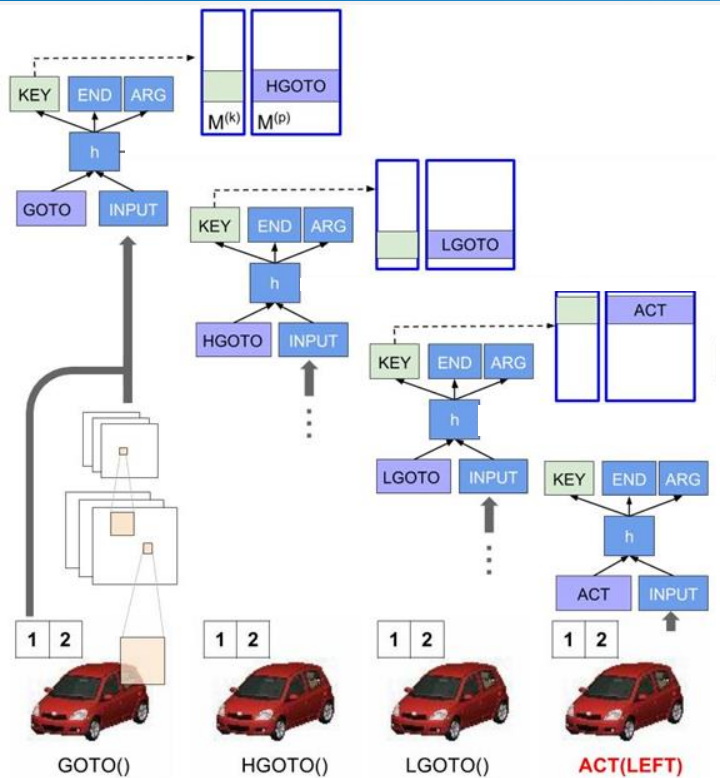
Environment observation

Input program

GOTO 1 2

**Generated commands**

GOTO 1 2

GOTO()

General Artificial Intelligence

GOTO()                    HGOTO()

GOTO()  HGOTO()  LGOTO()

General Artificial Intelligence

General Artificial Intelligence

GOTO()    HGOTO()    LGOTO()    **ACT(LEFT)**    LGOTO()    **ACT(LEFT)**    GOTO()    VGOTO()    DGOTO()    **ACT(DOWN)**

# Experiments

# Data Efficiency – Sorting

Seq2Seq LSTM and NPI used the same number of layers
and hidden units.

Trained on length 20 arrays of single-digit numbers.

NPI benefits from mining multiple subprogram examples per sorting instance, and additional parameters of the program memory.



Sorting per-sequence accuracy vs. # training examples

# Training examples

Seq2Seq      NPI

# Generalization – Sorting

For each length, we provided 64 example bubble sort traces, for a total of 1,216 examples.

Then, we evaluated whether the network can learn to sort arrays beyond length 20



**Sorting per-sequence accuracy vs sequence length**

Training sequence lengths

Sequence length

Seq2Seq    NPI

General Artificial Intelligence

# Generalization – Addition

Example problem: 90 + 160 = 250, we could represent the sequence as:

90X160X250

and solve addition via sequence prediction, e.g. "Learning to Execute" paper.

# Generalization – Addition problems

To make it easier, we can reverse and stack the inputs. (s2s-stack)

Even easier version: computation is entirely local. (s2s-easy)

output:  XXXX<span style="color:red">250</span>
input 1: <span style="color:blue">090</span>XXXX
input 2: <span style="color:green">061</span>XXXX

output:  <span style="color:red">052</span>
input 1: <span style="color:blue">090</span>
input 2: <span style="color:green">061</span>

# Generalization – Addition problems



Addition generalization: NPI vs Seq2Seq

Trained on problem sizes 1,..,20 digits.

# Generalization – Addition problems

s2s-stack



Addition generalization: NPI vs Seq2Seq

- NPI@32 per-sequence
- S2S-stack@32 per-character
- S2S-stack@512 per-character
- S2S-easy@32 per-sequence
- S2S-easy@64 per-sequence

*Test sequence length*

Google DeepMind

General Artificial Intelligence

# Generalization – Addition problems



Addition generalization: NPI vs Seq2Seq

s2s-easy

# Generalization – Addition problems



Addition generalization: NPI vs Seq2Seq

NPI

- NPI@32 per-sequence
- S2S-stack@32 per-character
- S2S-stack@512 per-character
- S2S-easy@32 per-sequence
- S2S-easy@64 per-sequence

Test sequence length

General Artificial Intelligence

# Multi-task NPI – Core is shared across all programs

| Program | Descriptions | Calls |
|---------|--------------|-------|
| ADD | Perform multi-digit addition | ADD1, LSHIFT |
| ADD1 | Perform single-digit addition | ACT, CARRY |
| CARRY | Mark a 1 in the carry row one unit left | ACT |
| LSHIFT | Shift a specified pointer one step left | ACT |
| RSHIFT | Shift a specified pointer one step right | ACT |
| ACT | Move a pointer or write to the scratch pad | - |
| BUBBLESORT | Perform bubble sort (ascending order) | BUBBLE, RESET |
| BUBBLE | Perform one sweep of pointers left to right | ACT, BSTEP |
| RESET | Move both pointers all the way left | LSHIFT |
| BSTEP | Conditionally swap and advance pointers | COMPSWAP, RSHIFT |
| COMPSWAP | Conditionally swap two elements | ACT |
| LSHIFT | Shift a specified pointer one step left | ACT |
| RSHIFT | Shift a specified pointer one step right | ACT |
| ACT | Swap two values at pointer locations or move a pointer | - |
| GOTO | Change 3D car pose to match the target | HGOTO, VGOTO |
| HGOTO | Move horizontally to the target angle | LGOTO, RGOTO |
| LGOTO | Move left to match the target angle | ACT |
| RGOTO | Move right to match the target angle | ACT |
| VGOTO | Move vertically to the target elevation | UGOTO, DGOTO |
| UGOTO | Move up to match the target elevation | ACT |
| DGOTO | Move down to match the target elevation | ACT |
| ACT | Move camera 15° up, down, left or right | - |
| RJMP | Move all pointers to the rightmost posiiton | RSHIFT |
| MAX | Find maximum element of an array | BUBBLESORT,RJMP |

# Learning new programs with a fixed NPI core

Toy example: Maximum-finding in an array.

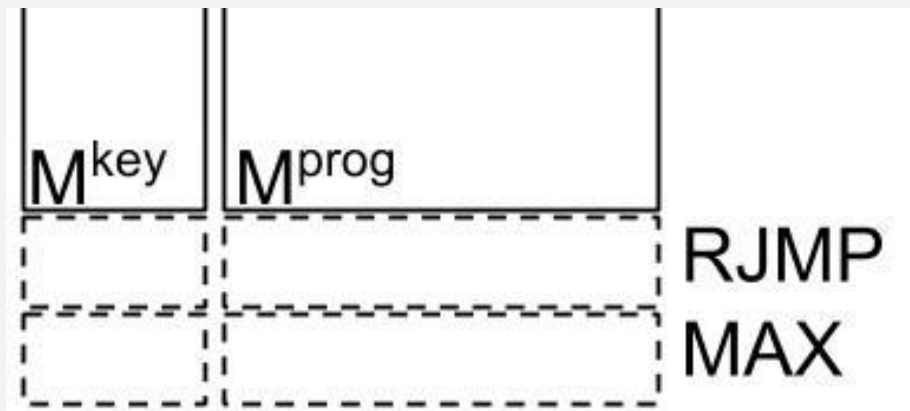Simple (not optimal) way: Call BUBBLESORT and then take the right-most element of the array. Two new programs:

**RJMP**: Move all pointers to the rightmost position in the array by repeatedly calling RSHIFT program.
**MAX**: Call BUBBLESORT and then RJMP

Expand program memory by adding 2 slots. Randomly initialize, then learn by backpropagation with the NPI core and all other parameters fixed.

# Learning new programs with a fixed NPI core



Protocol:
1. Randomly initialize new program vectors in memory
2. Freeze core and other program vectors
3. Backpropagate gradients to new program vectors

# Quantitative Results

| Task | Single | Multi | + Max |
|---|---|---|---|
| Addition | 100.0 | 97.0 | 97.0 |
| Sorting | 100.0 | 100.0 | 100.0 |
| Canon. seen car | 89.5 | 91.4 | 91.4 |
| Canon. unseen | 88.7 | 89.9 | 89.9 |
| Maximum | - | - | 100.0 |

- Per-sequence % accuracy.
- + Max indicates performance after addition of MAX program to memory.
- "unseen" uses a test set with disjoint car models from the training set.

# Conclusions & Next Steps

- A single NPI can learn multiple programs in dissimilar environments with different affordances.
- NPI sorting and addition programs exhibit strong generalization compared to baseline Seq2Seq models.
- A trained NPI with a fixed core can continue to learn new programs without forgetting already learned programs.

- **Next steps**: reduce supervision, scale up #programs, integrate new perception modules and affordances.

# Related work

Too much to cover in 20 minutes!
- Sigma-Pi Units (Rumelhart, 1986): activations of one network become the weights of a second network. Slowly changing network learns to control rapidly-changing network (Schmidhuber, 1992).
- Hierarchical RL (Sutton 1999, Dietterich 2000, Andre and Sutton 2001).
- Recent extensions of Seq2Seq: NTM, Pointer Networks, Memory Networks, Stack/Queue/Dequeue-augmented recurrent networks.
- Several other ICLR'16 papers on neural program induction. Main difference is that NPI explicitly incorporates compositional program structure.
- Recent models of prefrontal cognitive control (Donnarumma 2015).
- Learning to Execute (Zaremba 2014)

# Thanks!

# NPI single time step computation

Traces: $\quad \xi_t^{inp} : \{e_t, i_t, a_t\} \qquad \xi_t^{out} : \{i_{t+1}, a_{t+1}, r_t\}$

LSTM core input $\qquad s_t = f_{enc}(e_t, a_t)$

LSTM output $\qquad h_t = f_{lstm}(s_t, p_t, h_{t-1})$

pred. return prob $\qquad r_t = f_{end}(h_t)$

next program key $\qquad k_t = f_{prog}(h_t)$

next program args $\quad a_{t+1} = f_{arg}(h_t)$

# NPI single time step computation

Traces: $\xi_t^{inp} : \{e_t, i_t, a_t\}$ $\qquad$ $\xi_t^{out} : \{i_{t+1}, a_{t+1}, r_t\}$

$$i^* = \underset{i=1..N}{\arg\max} (M_{i,:}^{\mathbf{key}})^T k_t \;\;,\;\; p_{t+1} = M_{i^*,:}^{\mathbf{prog}}$$
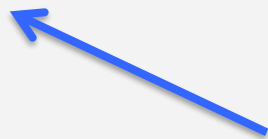
selected
prog inde

program
key memory

program
memory

x

# NPI single time step computation

**Traces:** $\xi_t^{inp} : \{e_t, i_t, a_t\}$     $\xi_t^{out} : \{i_{t+1}, a_{t+1}, r_t\}$

$$e_{t+1} \sim f_{env}(e_t, p_t, a_t)$$

Next environment observation; depends on selected program and arguments.
(Not controlled by NPI parameters)

# NPI learning

Traces: $\xi_t^{inp} : \{e_t, i_t, a_t\}$     $\xi_t^{out} : \{i_{t+1}, a_{t+1}, r_t\}$

Objective:

$$\theta^* = \arg\max_\theta \sum_{(\xi^{inp}, \xi^{out})} \log P(\xi^{out} | \xi^{inp}; \theta)$$

$$\log P(\xi_{out} | \xi_{inp}; \theta) = \sum_{t=1}^{T} \log P(\xi_t^{out} | \xi_1^{inp}, ..., \xi_t^{inp}; \theta)$$