

# Strong Generalization and Efficiency in Neural Programs

# Введение

## Две важных цели

Корректность и эффективность.

- Немногие алгоритмы смогли продемонстрировать сильное обобщение, т.е. корректность
- Эффективность либо было сложно показать, либо это не было главным направлением исследования

# Введение

## Neural Program Induction

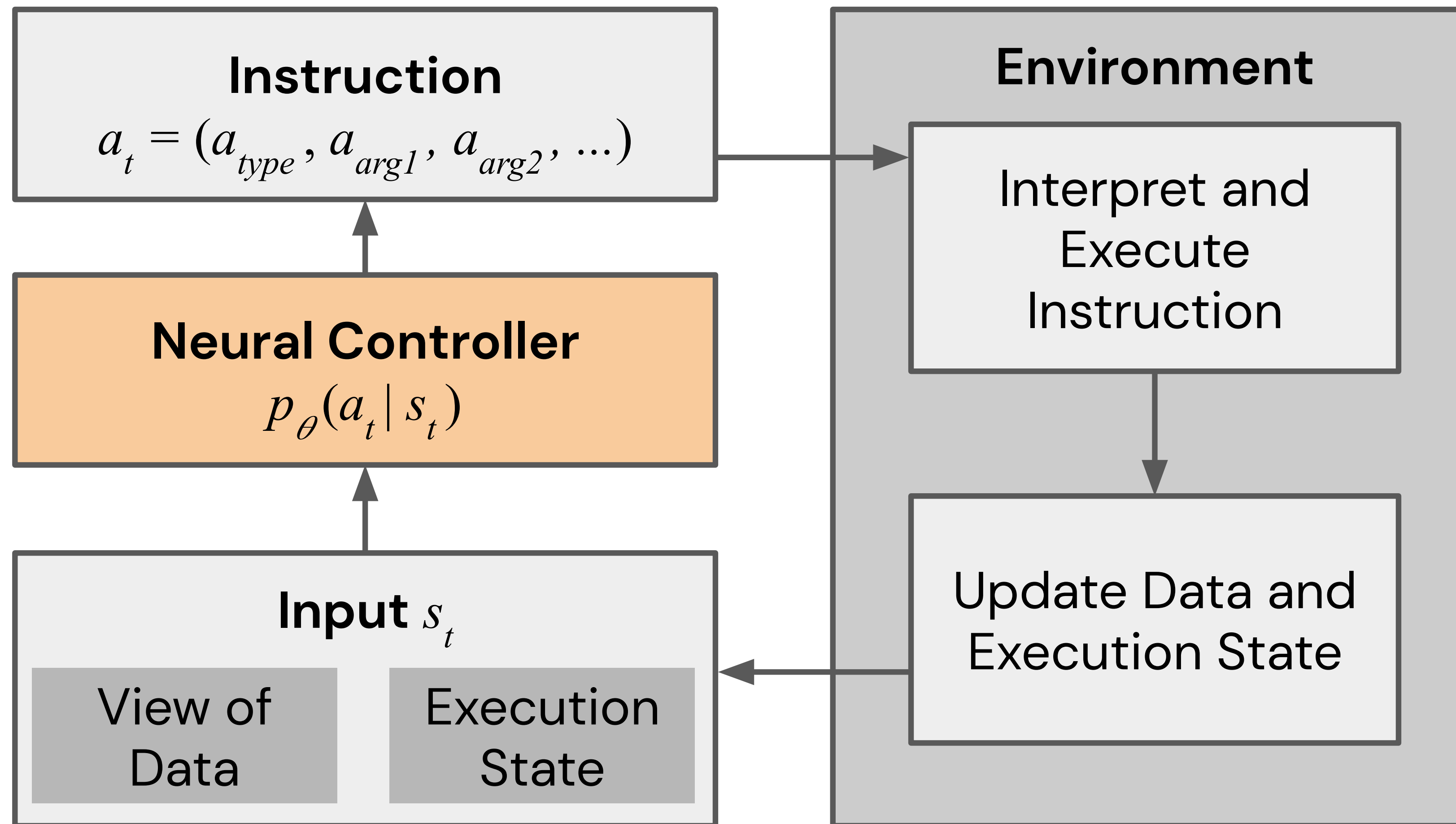
- Нейронный контроллер, который для состояния  $s$  возвращает инструкцию  $a = f(s)$
- Итерируется до критерия останова или максимального количества шагов
- Нейронная сеть параметризует policy распределения  $p(a | s)$

# Введение

Два важных результата статьи:

- Для эффективности и обобщения важно, какой вид у входа и выхода
- В то время, как обучение с учителем может подстроить алгоритм под учителя, RL может помочь превзойти его

# Описание метода



# Описание метода

Важные составляющие: входные данные, инструкции и модель.  
Для эффективности и обобщения используется константный размер входа.

В данном методе каждая инструкция имеет тип  $a_{type}$  и аргументы  $a_{arg1}, a_{arg2}, \dots$

Представление входных данных вместе с инструкциями формируют класс алгоритмов, который может представить нейронный контроллер.

# Описание метода

## Модель

Отображение  $a_t = f(s_t)$  определяется через  $p_\theta(a_t | s_t)$ :

$$p_\theta(a_t | s_t) = p_\theta(a_{type} | s_t) p_\theta(a_{args} | a_{type}, s_t)$$

$$p_\theta(a_{args} | a_{type}, s_t) = p_\theta(a_{arg1} | a_{type}, s_t) p_\theta(a_{arg2} | a_{arg1}, a_{type}, s_t)$$

В зависимости от того, имеет  $s_t$  фиксированный размер или переменный, мы используем или MLPs или GNNs.

# Обучение С учителем

Учитель генерирует  $(s_0, a_0, s_1, a_1, \dots)$

Минимизирует  $-\sum_t \log p_\theta(a_t | s_t)$

Быстро обучается, но не может превзойти учителя.



# Обучение RL

Награда за каждый шаг:  $r_t = -c$ .

Обновляем параметры  $\theta$  и  $\phi$  по направлению:

$[G_t - V_\phi(s_t)] \nabla_\theta \log p_\theta(a_t | s_t) + \mu \nabla_\phi [G_t - V_\phi(s_t)]^2$ , где

$$G_t = r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \gamma^n V_\phi(s_{t+n})$$

$\mu$  - скалярный вес

Добавляется лосс имитации  $-\lambda \log p_\theta(a'_t | s_t)$ , где  $a'_t = T(s_t)$

# Сортировка

## Переменный размер входа

Инструкция:  $\text{Swap}(i, j)$

$s_t$  содержит информацию о всем листе на каждом шаге. Устроен как полный граф из  $n$  вершин, поэтому размер входа переменный и используется GNNs.

Получаем сложность  $O(n^3)$

# Сортировка

## Константный размер входа

Обычные алгоритмы не смотрят на весь список.

1. Стоимость каждого шага становится константной
2. Алгоритм не зависит от длины последовательности

Будем работать с  $k$  индексами и включать во вход сравнение  $A[i]$  с  $A[i - 1]$  и  $A[i + 1]$ . Плюс для каждого  $j$ ,  $i$  будем знать сравнение  $A[i]$  с  $A[j]$

$$k = 4, v_1, v_3 = 0; v_2, v_4 = n - 1$$

Добавим инструкции: SwapWithNext( $i$ ), MoveVar( $i$ , +1/-1):  $v_i \leftarrow \min\{v_i, n - 1\}$  или  $\max\{v_i - 1, 0\}$ , AssignVar( $i$ ,  $j$ ):  $v_i \leftarrow v_j$

# Сортировка

## Результаты

[illegible]

# Повышение эффективности

## Функции позволяют разделять и властвовать

- $\text{FunctionCall}(id, l_1, \dots, l_p, o_1, \dots, o_p, r_1, \dots, r_p)$

$id$  - ID функции,  $l_i$  - ID локальной переменной,  $o_i$  - ID внешней переменной,  $r_i$  - ID внешней переменной, получающей возвращенное значение

- $\text{Return}(l'_1, \dots, l'_q)$

$l'_i$  - ID возвращаемой локальной переменной

$s_t$  содержит информацию, что и раньше, ID функции и кодировку прошлого состояния

Для задачи сортировки будем использовать 2 функции и  $\text{Swar}(i, j)$

При обучении с учителем получаем имитацию быстрой сортировки с  $O(n \log n)$

# Повышение эффективности

## RL превосходит учителя

Shaping награда: 
$$h(\mathbf{A}) = \sum_{i=low}^{high-1} \mathbb{I}[\mathbf{A}[i] \leq \mathbf{A}[i + 1]], \text{ когда RL from scratch}$$

Общая награда:  $r_t = h(\mathbf{A}_{t+1}) - h(\mathbf{A}_t) - c$

Обучение быстрее, но есть риск попасть в локальный минимум

Instance size	5	10	20	30	50	100	200	500	1000	10000
Bubble sort	13.5	68.0	293.6	677.8	1,874.7	7,527.0	30,046.8	187,506.2	750,519.3	-
Insertion sort	13.7	53.4	208.6	475.3	1,275.7	5,077.5	20,134.3	125,136.8	501,051.7	-
RL from scratch	8.7	49.7	252.0	617.5	1,763.7	7,320.5	29,623.0	186,305.1	748,784.5	-
RL + imitation	<b>8.2</b>	<b>44.3</b>	<b>190.7</b>	<b>446.6</b>	<b>1,228.5</b>	<b>4,981.2</b>	<b>19,939.2</b>	<b>124,643.0</b>	<b>500,057.0</b>	-
Quick sort	27.4*	87.5*	241.8	399.9	788.8	1,840.2	4,217.1	12,519.4	27,633.5	368,338.6
RL from scratch	<b>13.5*</b>	<b>56.0*</b>	<b>220.6</b>	495.5	1,315.6	5,166.5	20,321.4	125,622.0	502,035.5	-
RL + imitation	24.8*	79.9*	223.7	<b>377.2</b>	<b>728.3</b>	<b>1,717.4</b>	<b>3,914.4</b>	<b>11,578.2</b>	<b>25,635.0</b>	<b>338,947.9</b>

# Общность

## Поиск в упорядоченном списке

Используется тот же входной интерфейс, MoveVar, AssignVar

Добавляется:

- AssignMidVar(l, j, k):  $v_i \leftarrow \lfloor (v_j + v_k)/2 \rfloor$
- Found(i): элемент q найден по индексу  $v_i$
- NotFound(): элемент q не найден в **A**

Эта настройка определяет семейство контроллеров, которое включает в себя как алгоритм линейного сканирования  $O(n)$ , так и алгоритм двоичного поиска  $O(\log n)$ .

# Общность

## Поиск в упорядоченном списке

Instance Size	5	10	20	30	50	100	200	500	1000
Linear Search	2.2	4.6	9.7	14.2	23.3	50.0	106.7	220.7	446.6
Binary Search	2.6	3.9	6.3	7.5	9.2	11.7	<b>14.6</b>	<b>19.1</b>	<b>21.8</b>
RL from scratch	<b>1.3</b>	<b>2.5</b>	4.3	5.9	9.9	15.5	29.7	85.7	168.6
RL + imitation	1.4	2.6	<b>4.1</b>	<b>4.7</b>	<b>7.0</b>	<b>10.6</b>	18.4	41.5	85.8



# Общность

## 0/1 задача о рюкзаке

$$w_i, v_i \sim U[0,1], W = \frac{1}{2} \sum_i w_i$$

Инструкции: Pop(i), Put(i)

- MoveVar(+1/-1):  $i \leftarrow i + 1$
- Knapsack() - рекурсивно вызывает функцию knapsack
- Return() - возвращает текущий уровень рюкзака

Это семейство алгоритмов включает DFS

$r_t =$  (лучшая ценность после  $a_t$ ) — (лучшая ценность до  $a_t$ )

# Общность

## 0/1 задача о рюкзаке

Budget	Size	2	4	6	8	10	20	40	80	160	320	640
20x size	DFS	0.45	1.25	1.82	2.47	3.14	5.94	11.10	21.67	42.11	82.01	163.42
	RL	0.45	<b>1.26</b>	<b>2.03</b>	<b>2.83</b>	<b>3.51</b>	<b>6.51</b>	<b>11.78</b>	<b>22.40</b>	<b>43.27</b>	<b>83.07</b>	<b>164.56</b>
100x size	DFS	0.45	1.26	<b>2.04</b>	2.81	3.43	6.43	11.75	22.39	42.95	82.82	164.22
	RL	0.45	1.26	2.03	<b>2.86</b>	<b>3.70</b>	<b>6.98</b>	<b>12.41</b>	<b>23.15</b>	<b>44.16</b>	<b>84.05</b>	<b>165.60</b>

# Вопросы

1. Распишите, как выглядит policy с объяснением.
2. Как выглядит reward вместе с shaping reward и зачем его добавлять?
3. Расскажите о приложении к задаче о рюкзаке.