

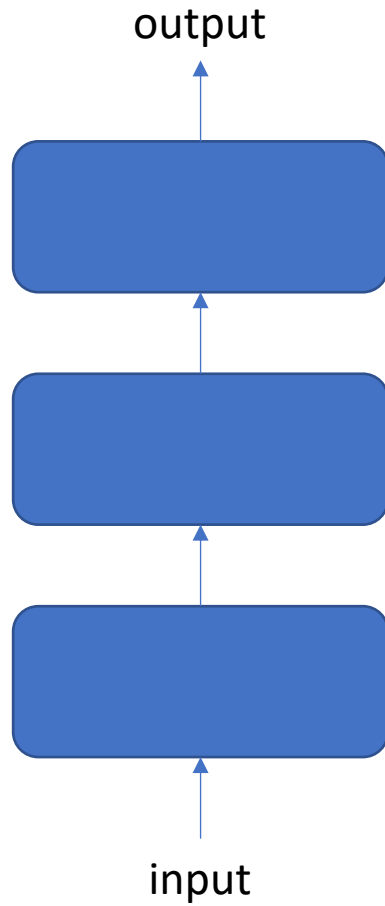
# Progressive Neural Networks

Чистяков Глеб, гр. 162

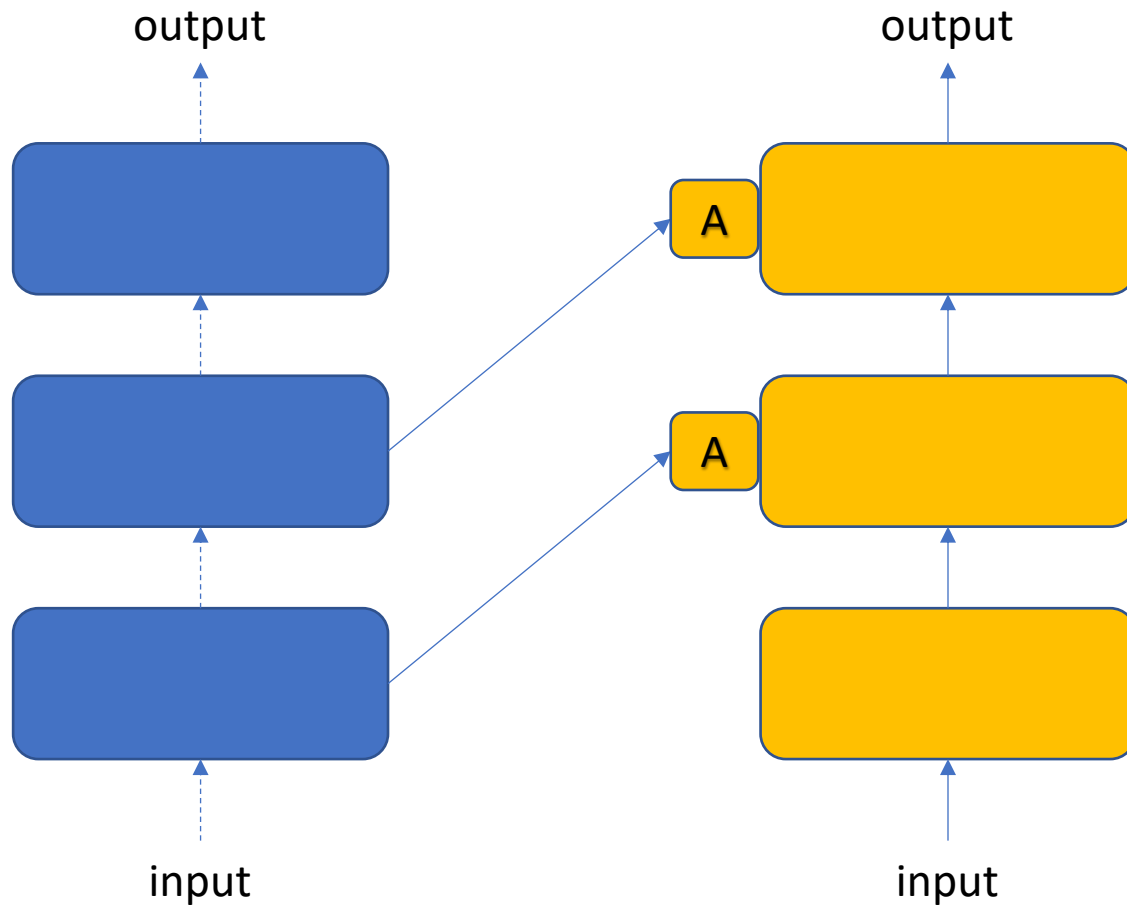
# Progressive Networks

- Способность включать предварительные знания на каждом уровне иерархии объектов
- Способность повторно использовать старые вычисления и изучать новые
- Невосприимчивость к катастрофическому забвению (catastrophic forgetting)

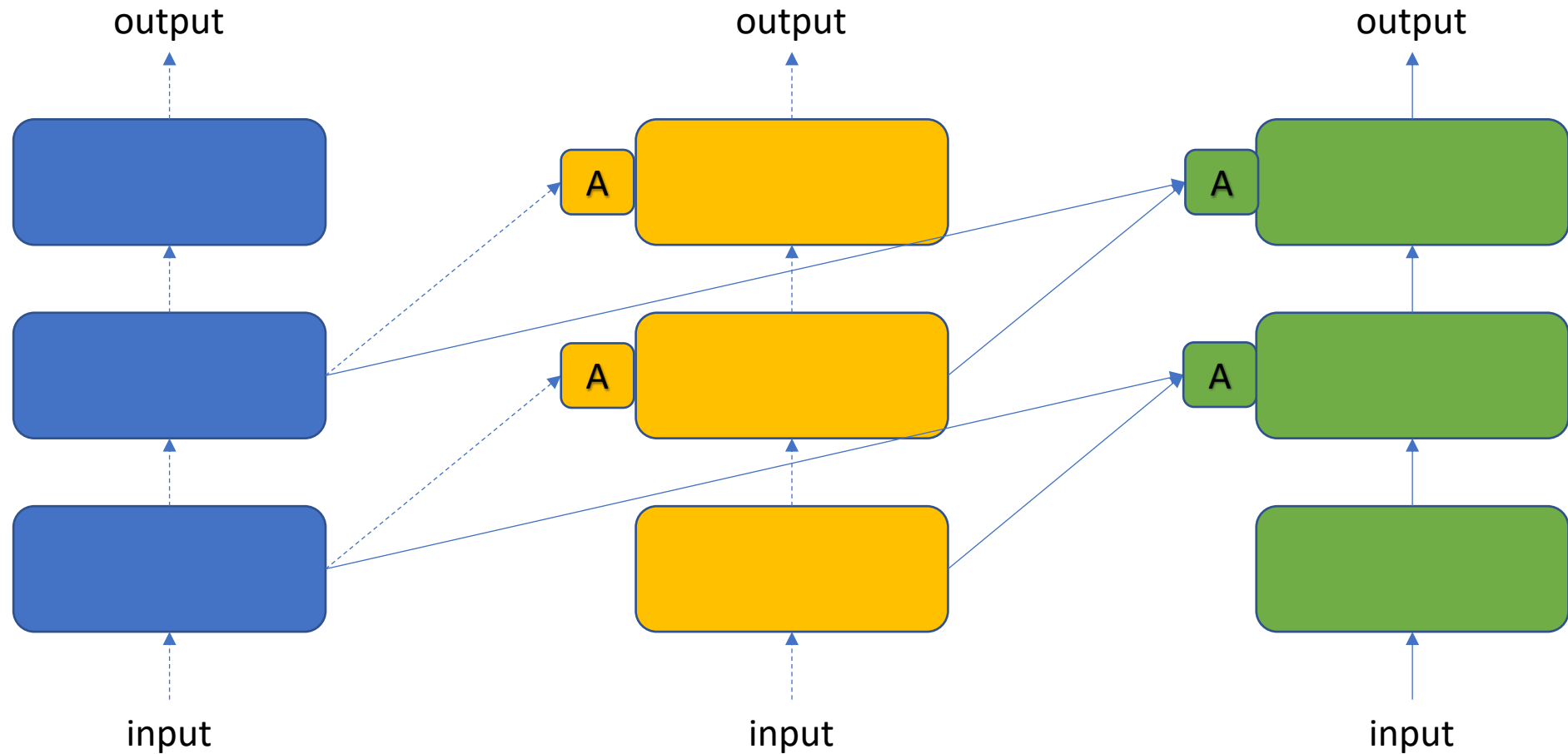
# Progressive Networks



# Progressive Networks



# Progressive Networks



# Hidden activations

$$h_i^{(k)} = f \left( W_i^{(k)} h_{i-1}^{(k)} + \sum_{j < k} U_i^{(k:j)} h_{i-1}^{(j)} \right)$$

$W_i^{(k)} \in \mathbb{R}^{n_i \times n_{i-1}}$  — матрица весов слоя  $i$  столбца  $k$

$U_i^{(k:j)} \in \mathbb{R}^{n_i \times n_j}$  — боковые соединения от слоя  $i - 1$  столбца  $j$  до слоя  $i - 1$  столбца  $k$

$f(x) = \max(0, x)$

# Adapters

$h_{i-1}^{(<k)} = [h_{i-1}^{(1)} \dots h_{i-1}^{(j)} \dots h_{i-1}^{(k-1)}]$  – вектор предыдущих параметров

$$h_i^{(k)} = \sigma \left( W_i^{(k)} h_{i-1}^{(k)} + U_i^{(k:j)} \sigma(V_i^{(k:j)} \alpha_{i-1}^{(<k)} h_{i-1}^{(<k)}) \right)$$

$V_i^{(k:j)} \in \mathbb{R}^{n_{i-1} \times n_{i-1}^{(<k)}}$  – матрица проекции

# Transfer Analysis

- Average Perturbation Sensitivity (APS)
- Average Fisher Sensitivity (AFS)



# Transfer Analysis

## Average Perturbation Sensitivity (APS)

$$\Lambda_i^{(k)} = \frac{1}{\sigma_i^{2(k)}} \text{ — точность шума, вводимого в слой } i \text{ столбца } k$$

$$\text{APS}(i, k) = \frac{\Lambda_i^{(k)}}{\sum_k \Lambda_i^{(k)}} \text{ — оценка чувствительности слоя } i \text{ столбца } k$$

# Transfer Analysis

## Average Fisher Sensitivity (AFS)

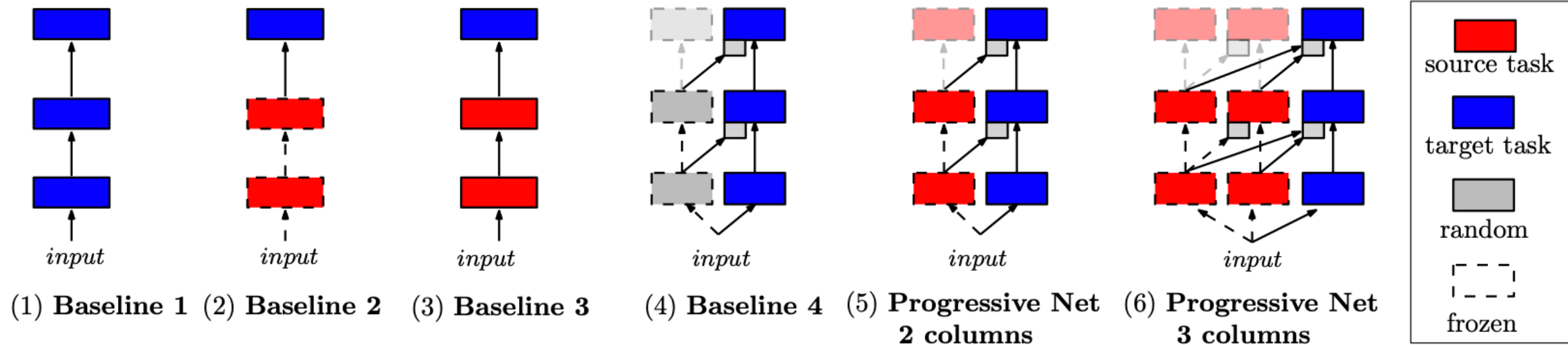
$$\hat{F}_i^{(k)} = \mathbb{E}_{\rho(s,a)} \left[ \frac{\partial \log \pi}{\partial \hat{h}_i^{(k)}} \frac{\partial \log \pi^T}{\partial \hat{h}_i^{(k)}} \right] \quad - \text{ матрица Фишера}$$

$$\pi^{(k)}(a \mid s) := h_L^{(k)}(s) \quad - \text{ Политика } k\text{-го столбца, принимающая состояние окружения и выдающая вероятность над действиями}$$

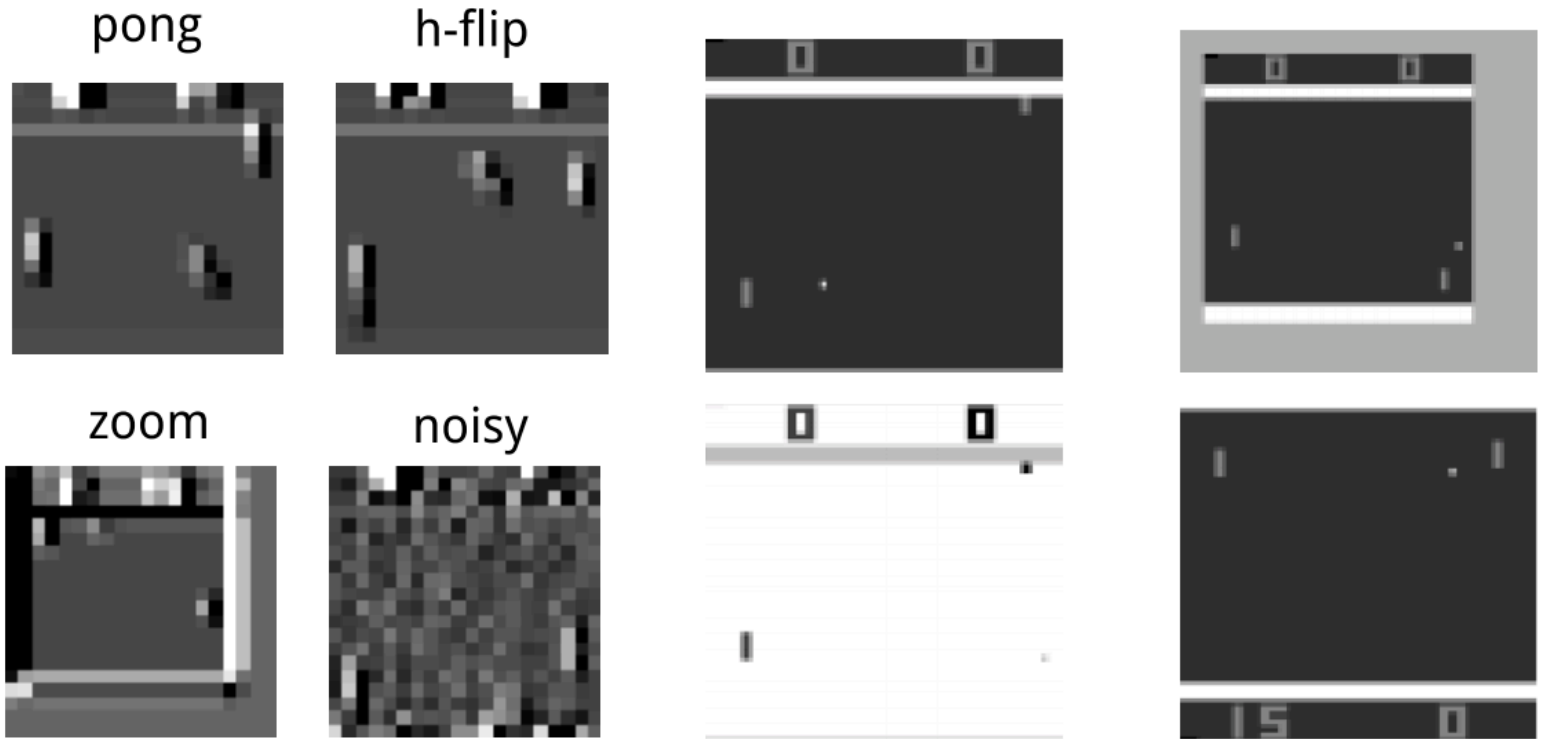
$$\text{AFS}(i, k, m) = \frac{\hat{F}_i^{(k)}(m, m)}{\sum_k \hat{F}_i^{(k)}(m, m)} \quad - \text{ средняя чувствительность Фишера признака } m \text{ слоя } i \text{ столбца } k$$

$$\text{AFS}(i, k) = \sum_m \text{AFS}(i, k, m) \quad - \text{ оценка чувствительности слоя } i \text{ столбца } k$$

# Experiments

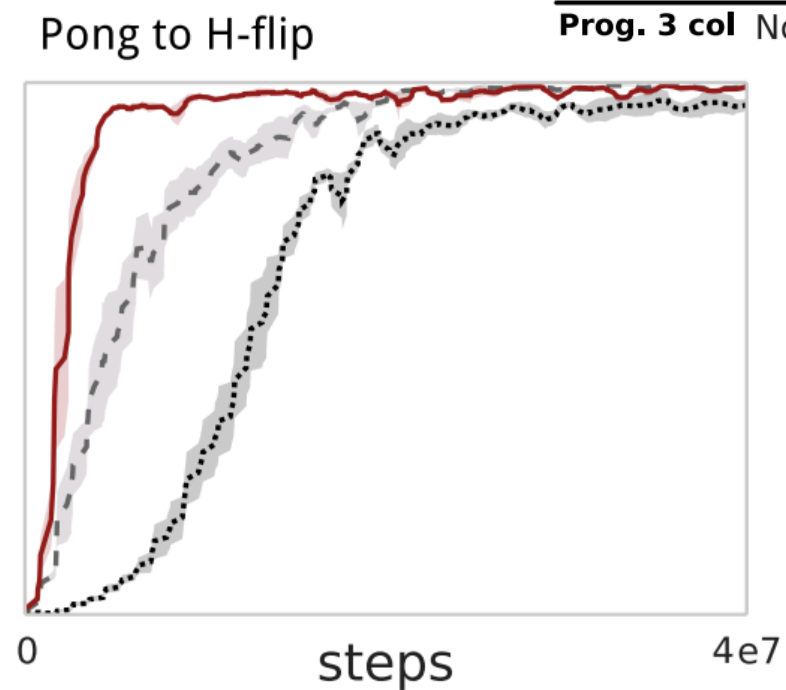
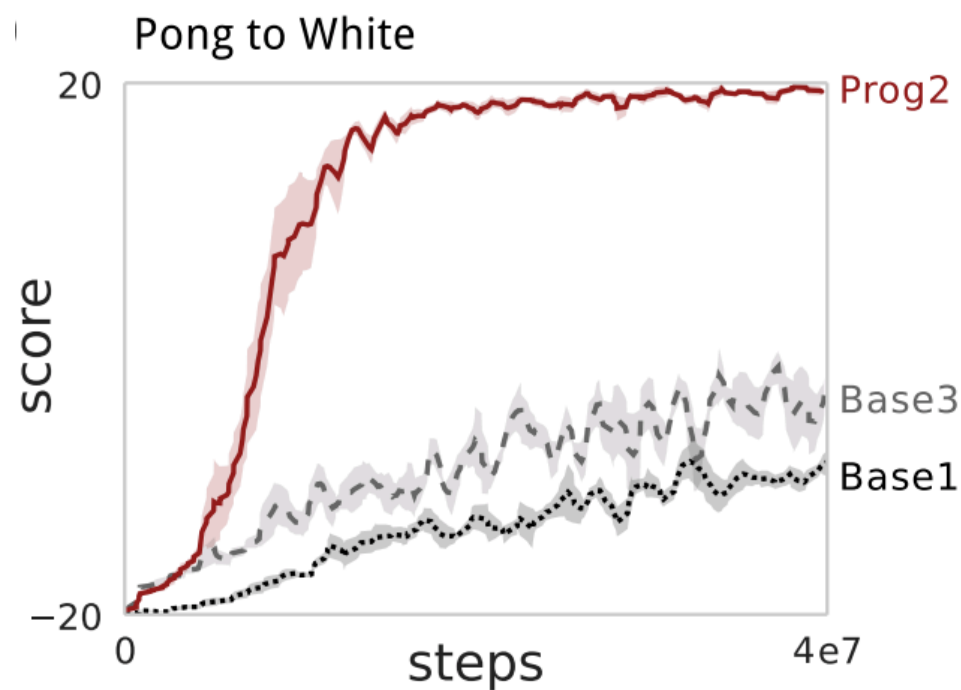


# Pong Soup



- Noisy (frozen Gaussian noise is added to the inputs)
- Black (black background)
- White (white background)
- Zoom (input is scaled by 75% and translated)
- V-flip (input is vertically flipped)
- H-flip (input is horizontally flipped)
- VH-flip (input is horizontally and vertically flipped)

# Pong Soup



		Pong	Black	H-flip	HV-flip	Noisy	V-flip	White	Zoom
<b>Baseline 2</b>	Pong	✖	1	0	0	0	0	1	0
	Noisy	1	1	0	✖	✖	1	1	0
	H-flip	0	0	✖	0	0	0	0	0
<b>Baseline 3</b>	Pong	✖	1	1	1	1	1	1	1
	Noisy	1	1	1	✖	1	1	1	1
	H-flip	1	1	✖	1	1	1	1	1
<b>Baseline 4</b>	Random	1	1	1	1	1	1	1	1
<b>Prog. 2 col</b>	Pong	✖	1	1	1	1	1	1	1
	Noisy	1	1	1	✖	1	1	1	1
	H-flip	1	1	✖	1	1	1	1	1
<b>Prog. 3 col</b>	Noisy + H-flip	1	✖	✖	1	✖	1	1	1

2

1

0

# Atari games

## Source games:

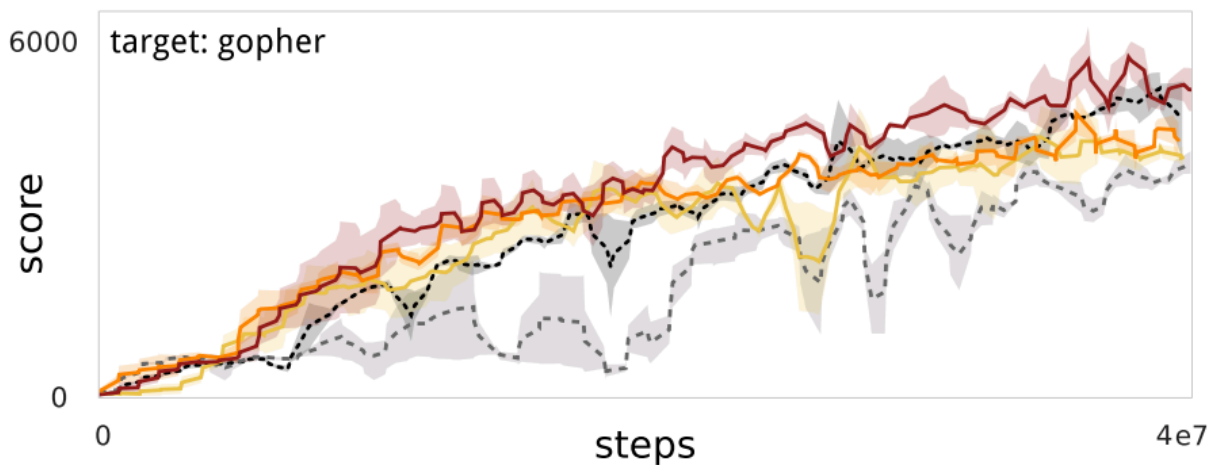
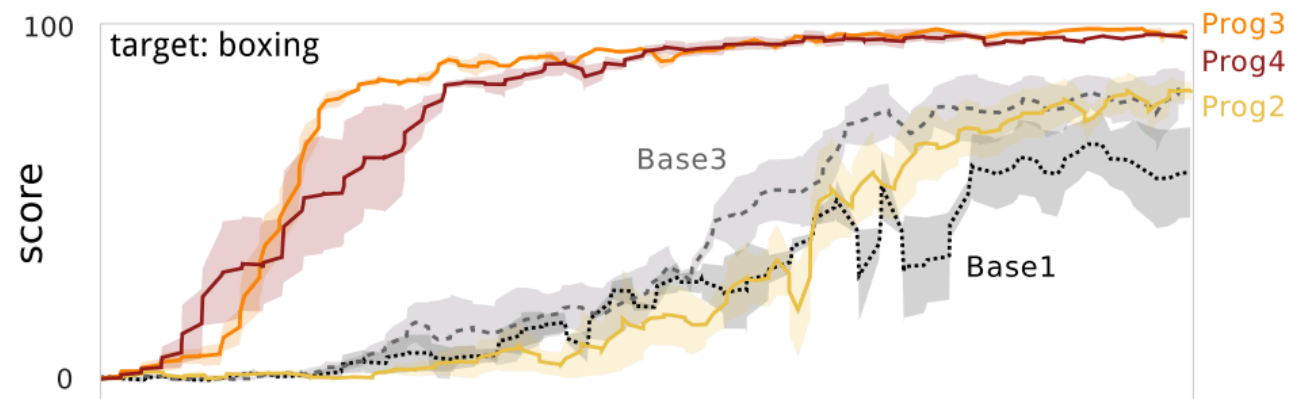
- Pong
- River Raid
- Seaquest

## Target games:

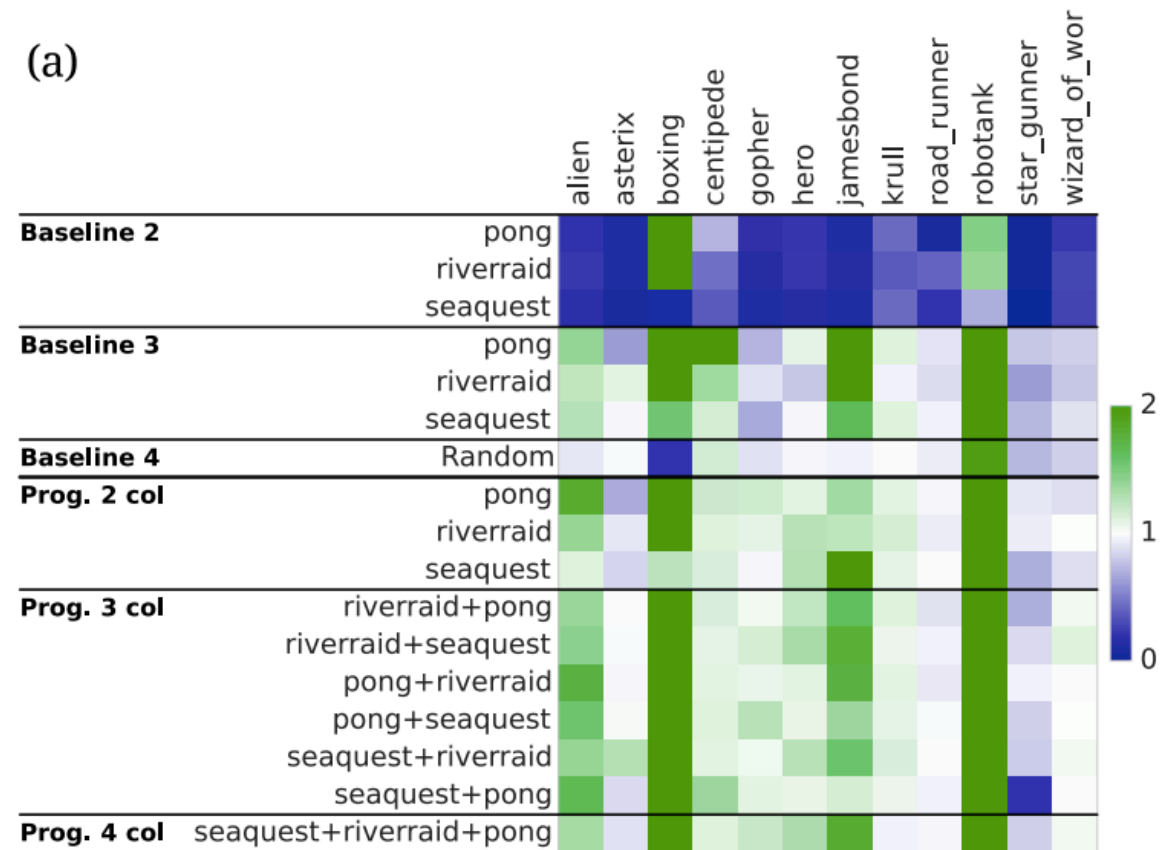
- Alien
- Asterix
- Boxing
- Centipede
- Gopher
- Hero
- James Bond
- Krull Robotank
- Road Runner
- Star Gunner
- Wizard of Wor



# Atari games

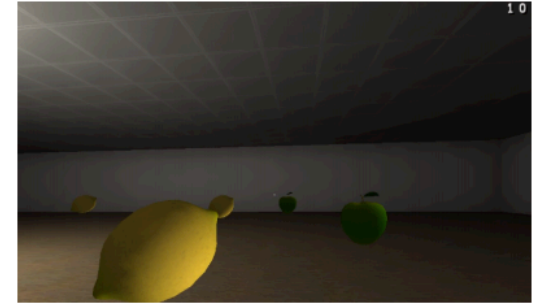
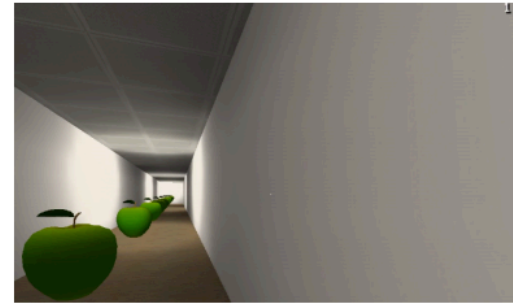


(a)



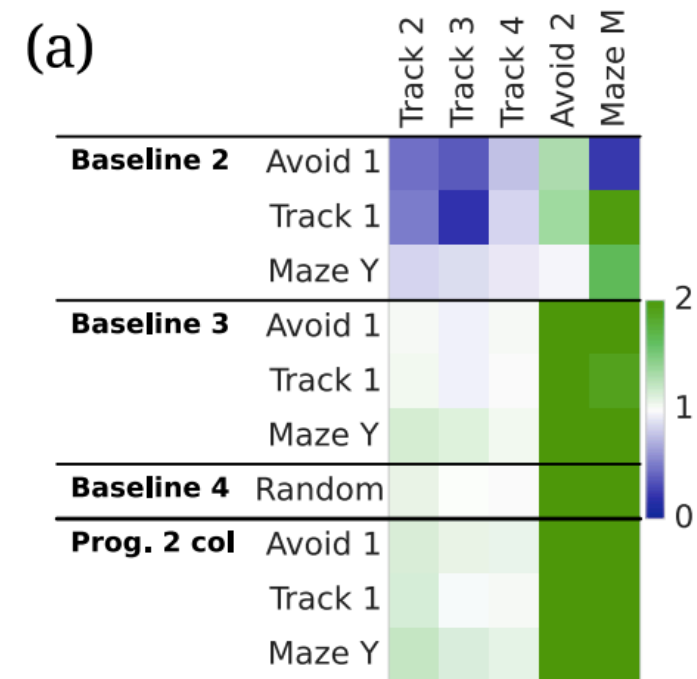
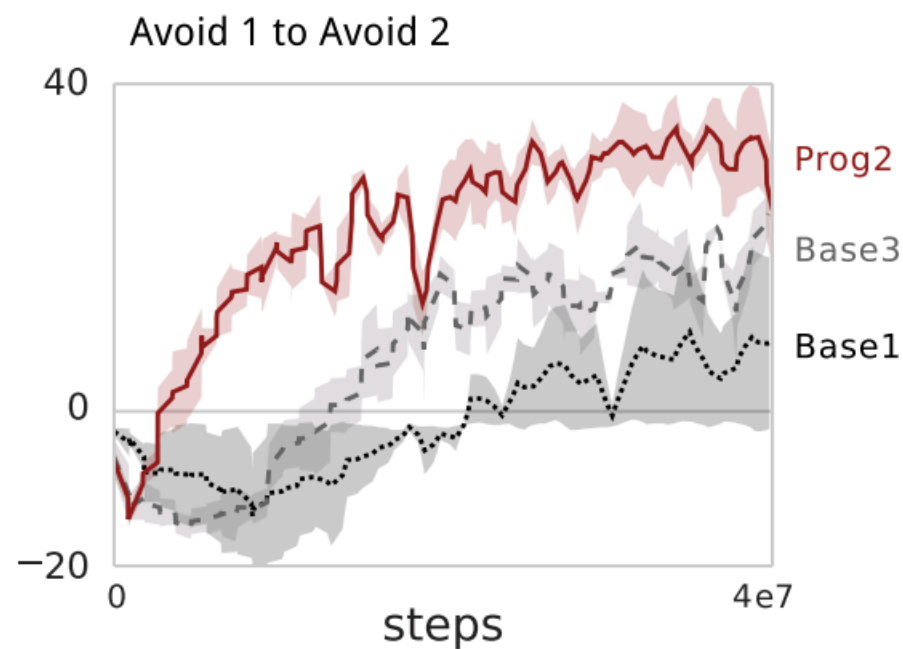
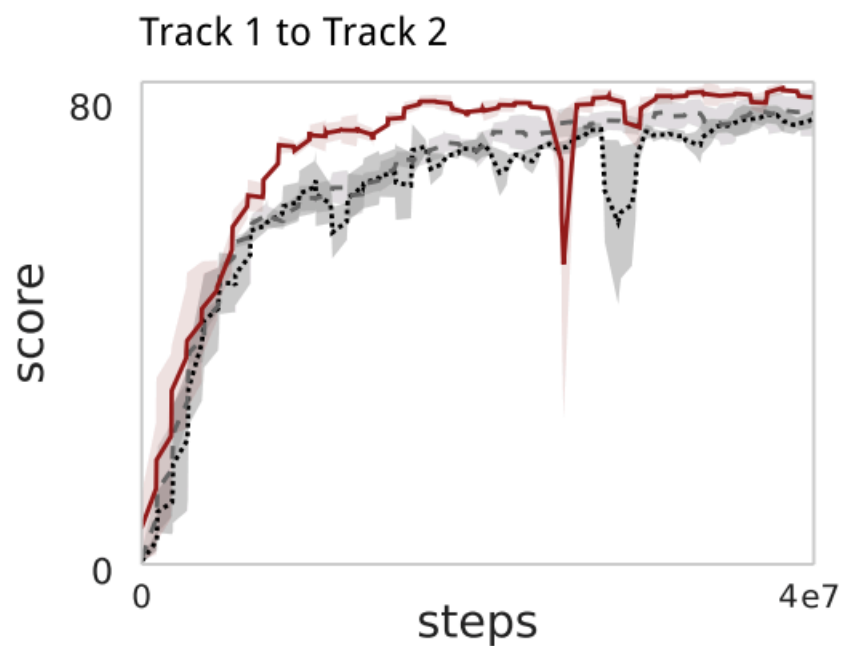
# Labyrinth games

- Seek Track 1: simple corridor with many apples
- Seek Track 2: U-shaped corridor with many strawberries
- Seek Track 3:  $\Omega$ -shaped, with 90o turns, with few apples
- Seek Track 4:  $\Omega$ -shaped, with 45o turns, with few apples
- Seek Avoid 1: large square room with apples and lemons
- Seek Avoid 2: large square room with apples and mushrooms
- Seek Maze M : M-shaped maze, with apples at dead-ends
- Seek Maze Y : Y-shaped maze, with apples at dead-ends





# Labyrinth games



# Experiments

	<b>Pong Soup</b>		<b>Atari</b>		<b>Labyrinth</b>	
	Mean (%)	Median (%)	Mean (%)	Median (%)	Mean (%)	Median (%)
Baseline 1	100	100	100	100	100	100
Baseline 2	35	7	41	21	88	85
Baseline 3	181	160	133	110	235	112
Baseline 4	134	131	96	95	185	108
Progressive 2 col	209	169	132	112	<b>491</b>	<b>115</b>
Progressive 3 col	<b>222</b>	<b>183</b>	140	111	—	—
Progressive 4 col	—	—	<b>141</b>	<b>116</b>	—	—

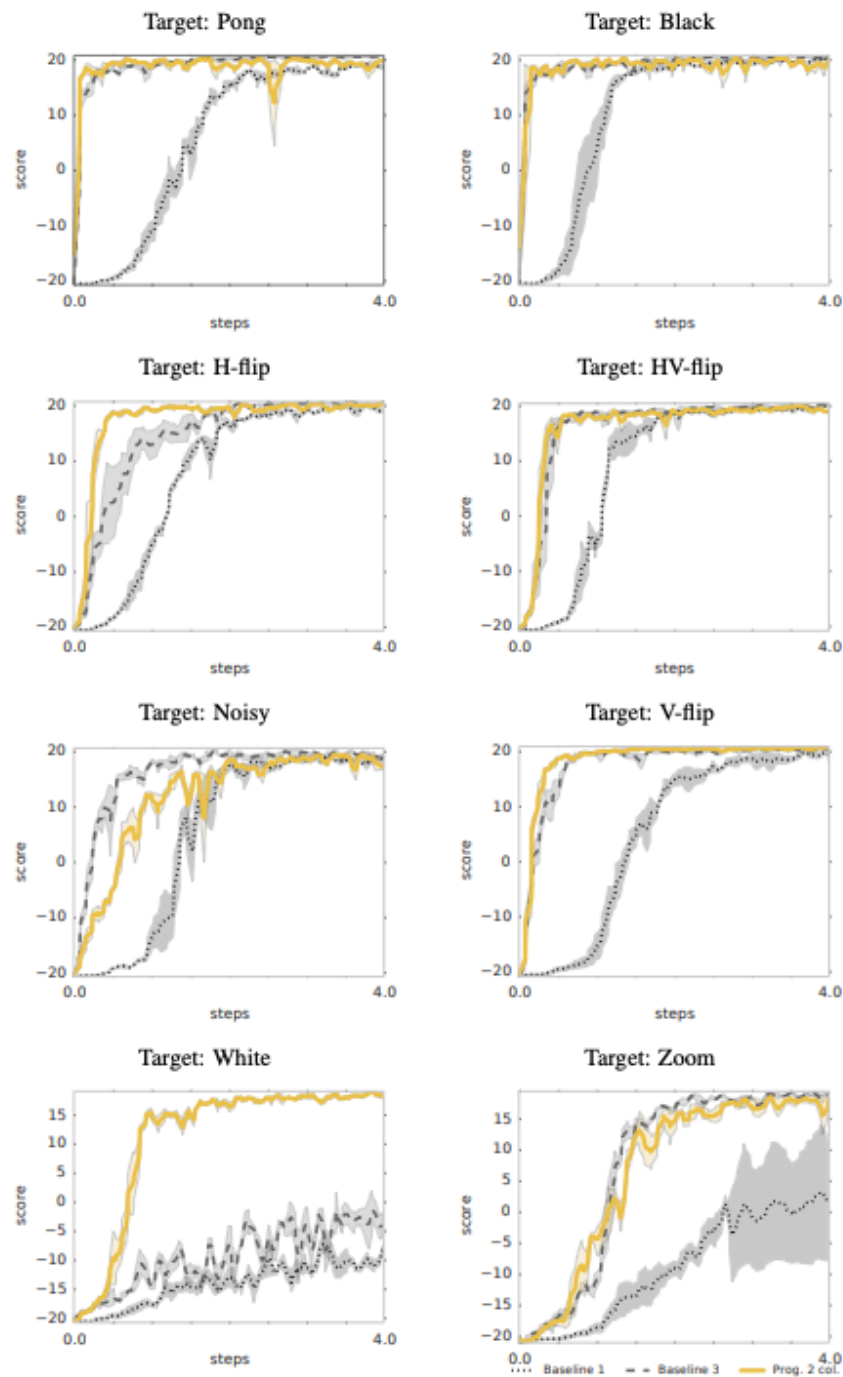


Figure 12: Training curves for transferring to 8 *target* games after learning standard Pong first.

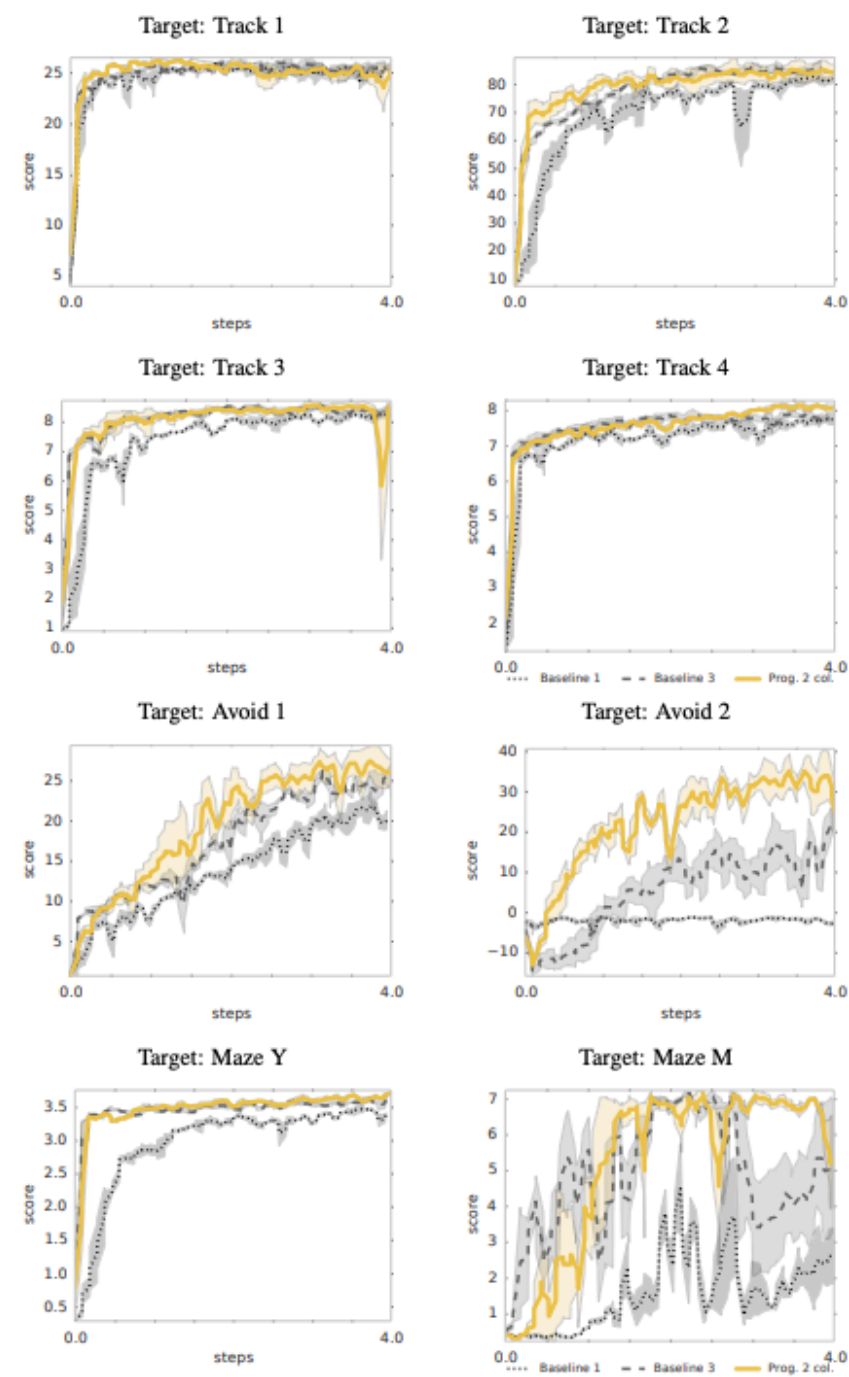


Figure 13: Training curves for transferring to 8 *target* games after learning Maze Y first.

# Quiz

- Какая идея у метода прогрессивных сетей? Написать формулу скрытой активации.
- Изобразить пример для сети из трех задач (три столбца). Описать интуицию метода.
- Что такое адаптеры? Написать формулу скрытого слоя адаптера.