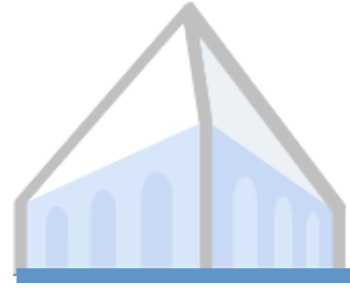# Modular multitask reinforcement learning with policy sketches
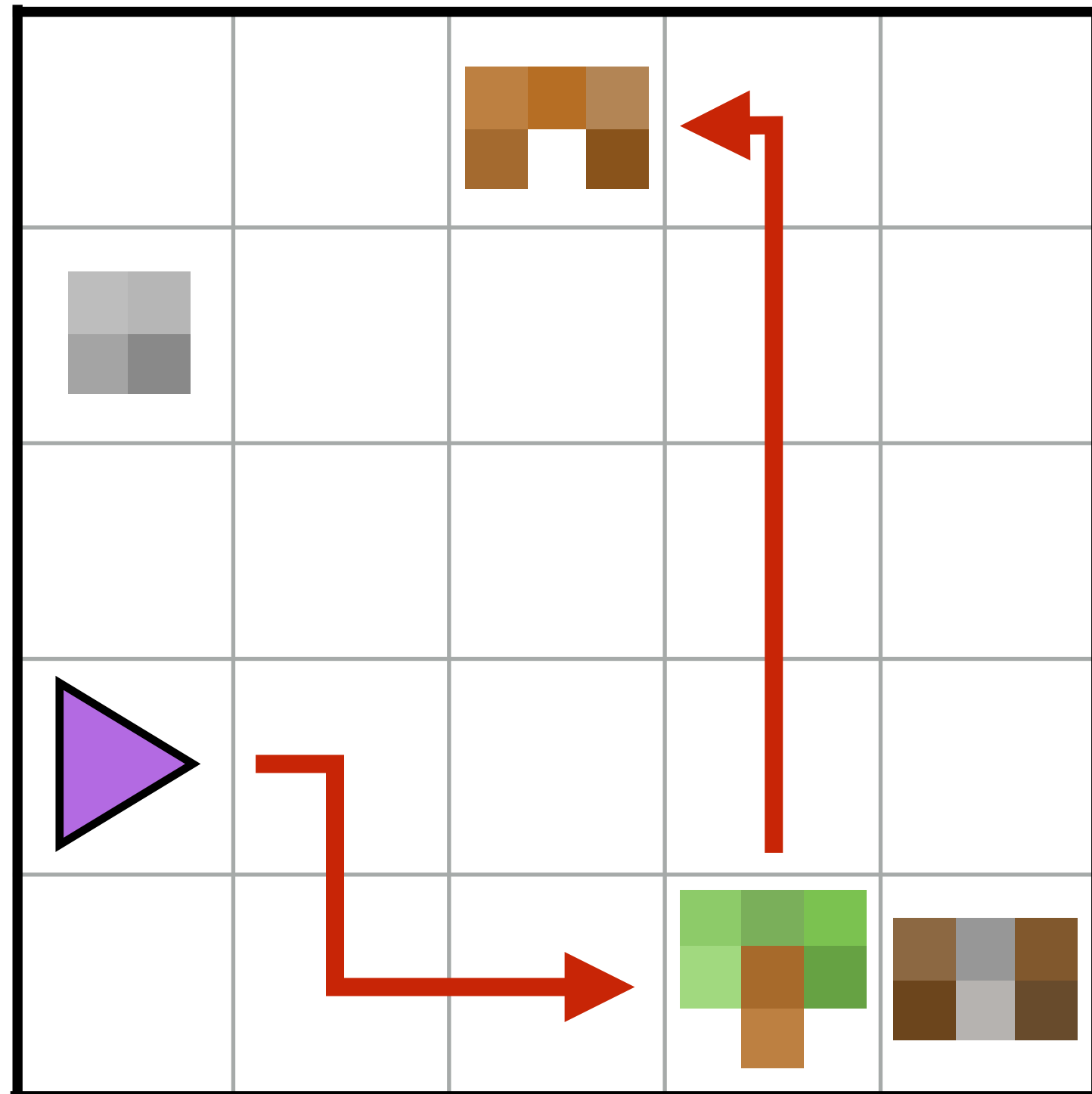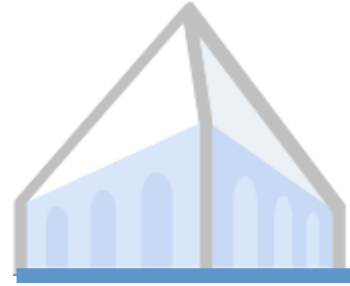


Jacob Andreas, Sergey Levine and Dan Klein

# The learning problem

**make planks**
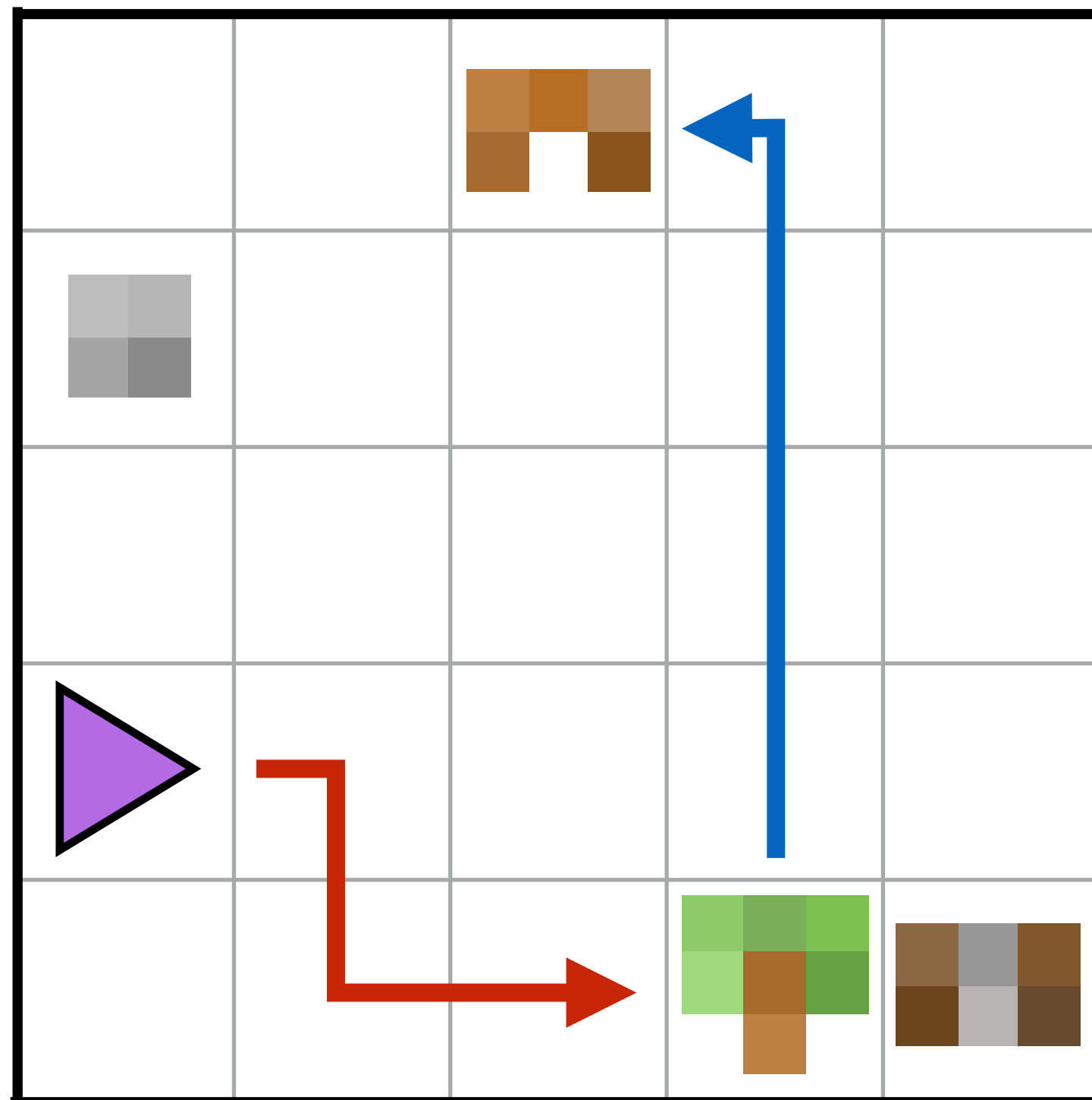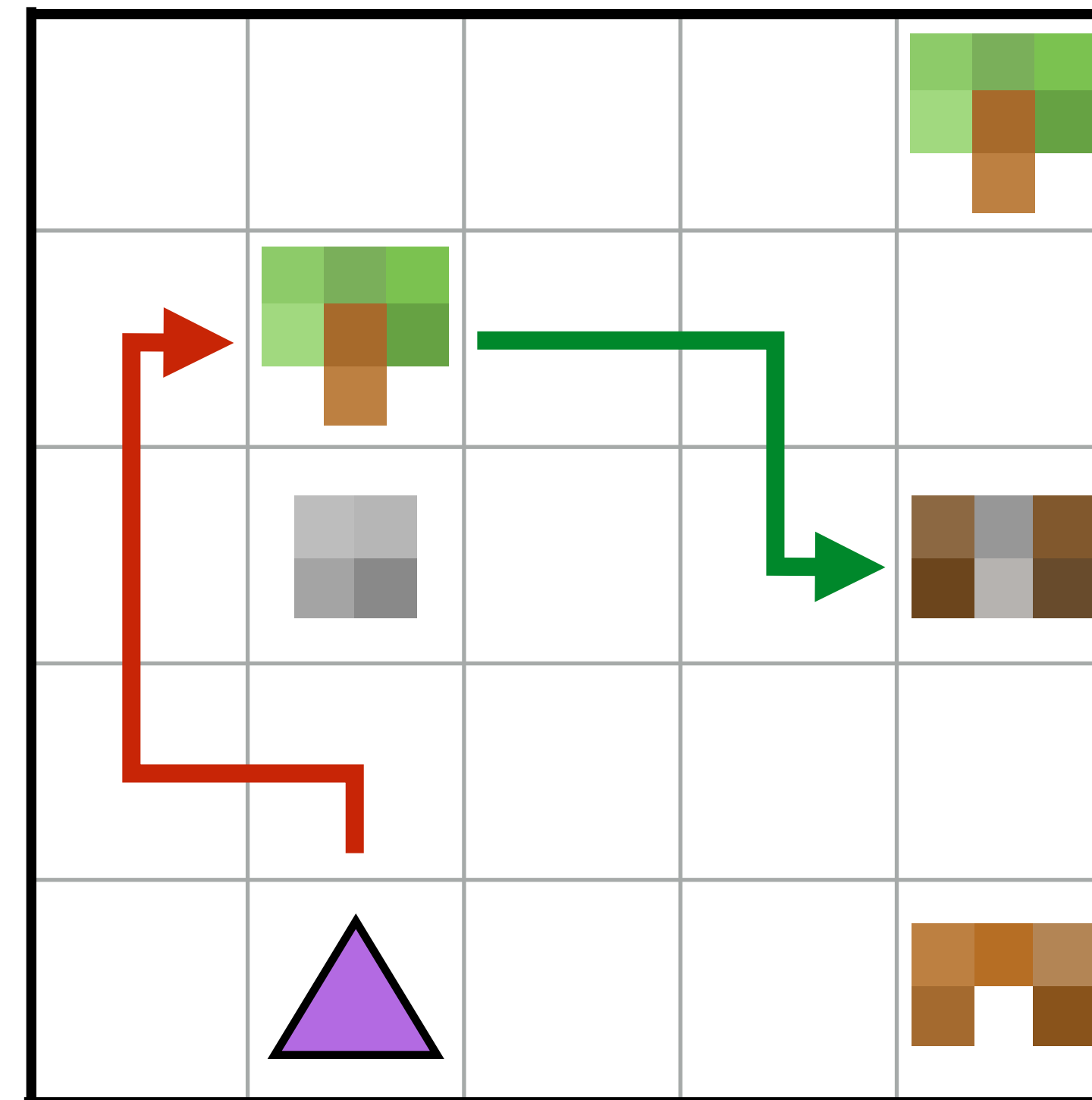
**make sticks**

# Learning from sketches

# The options framework

# The options framework

# Learning from intermediate rewards



[Kearns & Singh 02, Kulkarni et al. 16]

# Learning from demonstrations

[Stolle & Precup 02, Fox & Krishnan et al. 16]

# Learning from policy sketches
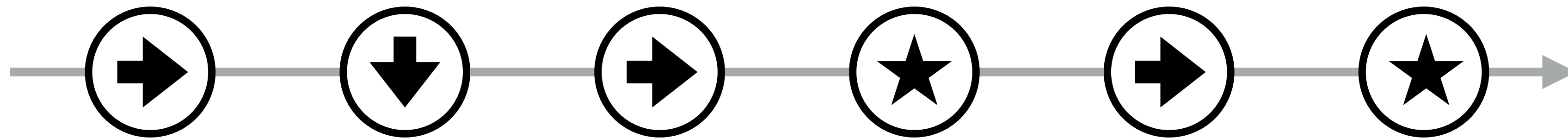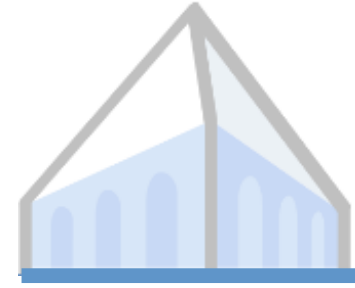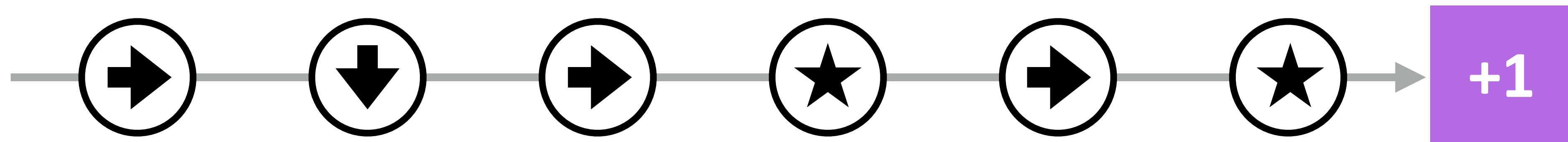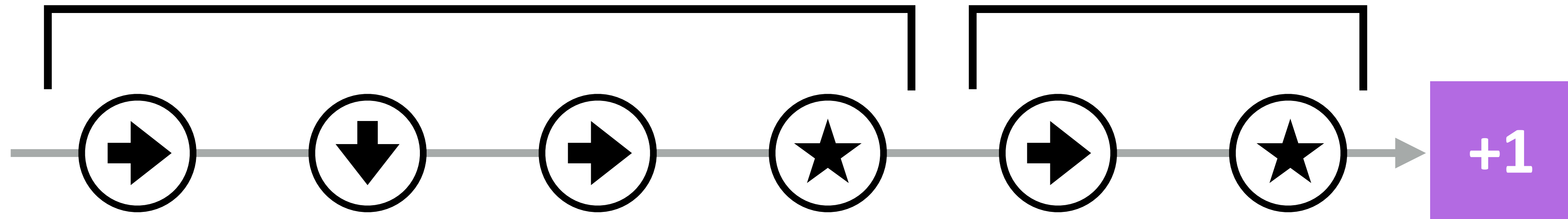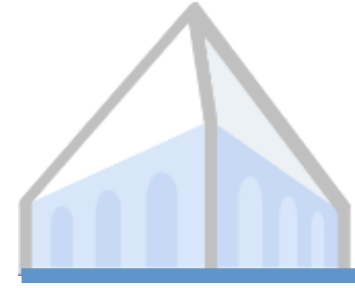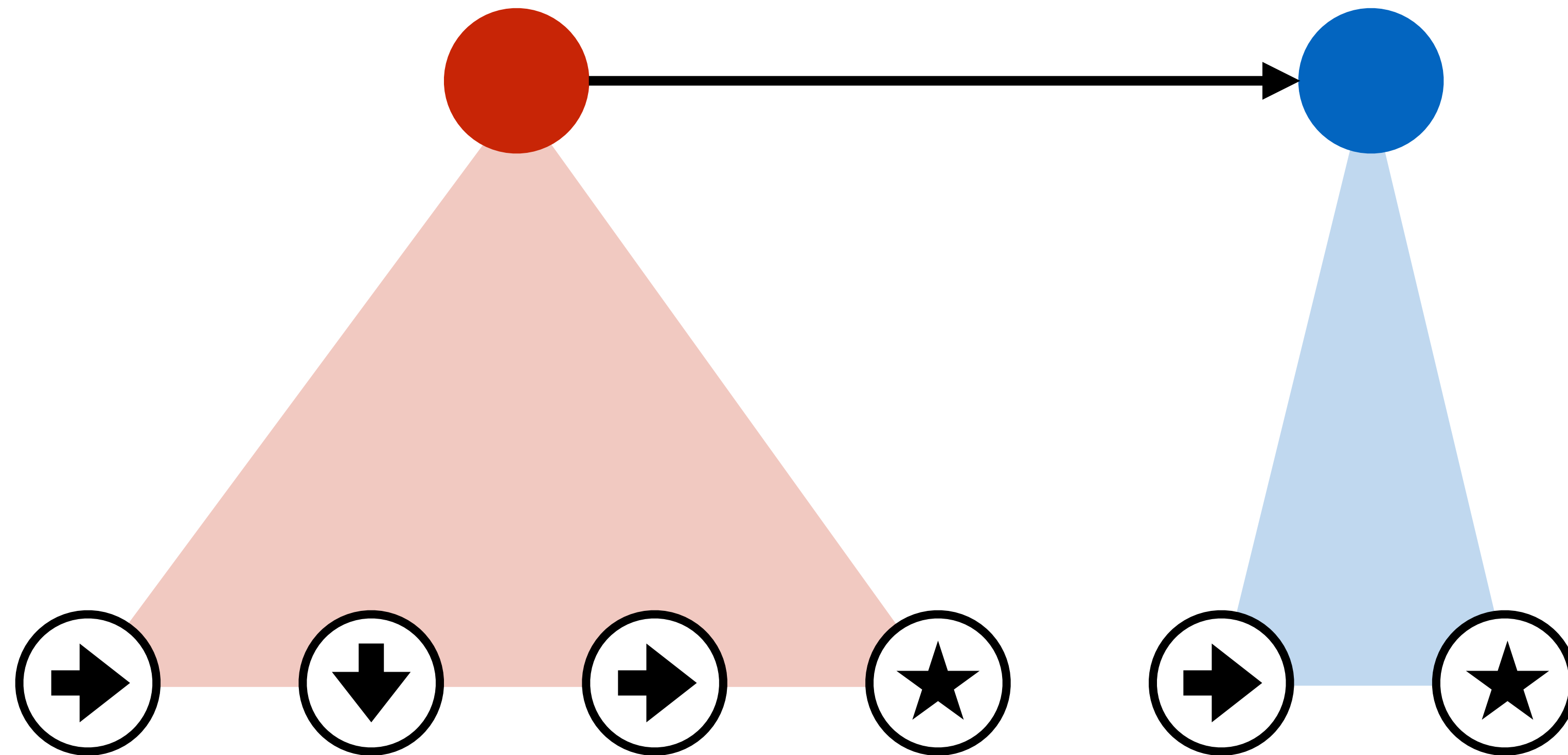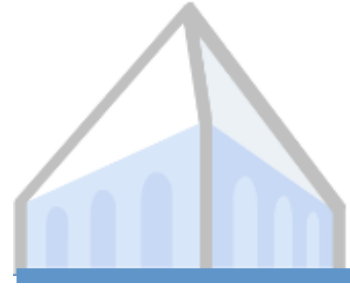
# Why sketches?

## Easy to collect

```
make plank    get wood    use tool
make stick    get wood    use work
make cloth    get grass   use fact
make rope     get grass   use tool
make bridge   get iron    get wood
make bed*     get wood    use tool
make axe*     get wood    use work
make shears   get wood    use work
get gold      get iron    get wood
get gem       get wood    use work
```

## Portable

# Learning from policy sketches

make planks

get wood

use saw

**make sticks**

**get wood**

**use axe**

# Learning from policy sketches

[e.g. Branavan et al. 09, Oh et al. 17, Hermann et al. 17]

get wood

use saw

get wood

use axe

get wood

use saw

$\pi_1$

$\pi_2$

get wood

use axe

$\pi_1$

$\pi_3$

get wood

use saw

$\pi_1$

$\pi_2$

get wood

use axe

$\pi_1$

$\pi_3$

get wood

$\pi_1$

$\pi_1$

get wood

???

$\pi_1$

get wood

# Action probabilities

$\pi_1$

get wood

# Policy search

action   state   reward   baseline

$$\sum_{tasks} \sum_{steps} \left( \nabla \log \pi(\blacktriangleright \mid \text{▦}) \right) (r_t - b)$$

$$\sum_{\text{tasks}} \sum_{\text{steps}} \left( \nabla \log \pi( \rightarrow | \blacksquare ) \right) (r_t - b)$$

get wood

$$\sum_{\text{tasks}} \sum_{\text{steps}} \left( \nabla \log \pi(\; \blacktriangleright \; | \; \square \;) \right) (r_t - b)$$

use axe

# Policy search

Reward

.40

$$\sum_{tasks} \sum_{steps} \left( \nabla \log \boxed{\text{SUBPOLICY}} \right) (r_t - b)$$

# Improving policy search

action   state   reward   baseline

$$\sum_{\text{tasks}} \sum_{\text{steps}} \left( \nabla \log \pi(\rightarrow | \text{▦}) \right) (r_t - b)$$

# Improving policy search

$$\left( \nabla \log \boxed{\texttt{use saw}} \right) \left( r_t - \boxed{\texttt{make planks}} \right)$$

$$\left( \nabla \log \boxed{\texttt{use saw}} \right) \left( r_t - \boxed{\texttt{make nails}} \right)$$

$$\left( \nabla \log \boxed{\texttt{use axe}} \right) \left( r_t - \boxed{\texttt{make planks}} \right)$$

$$\left( \nabla \log \boxed{\texttt{use axe}} \right) \left( r_t - \boxed{\texttt{make nails}} \right)$$

$$\left( \nabla \log \boxed{\texttt{get wood}} \right) \left( r_t - \boxed{\texttt{make planks}} \right)$$
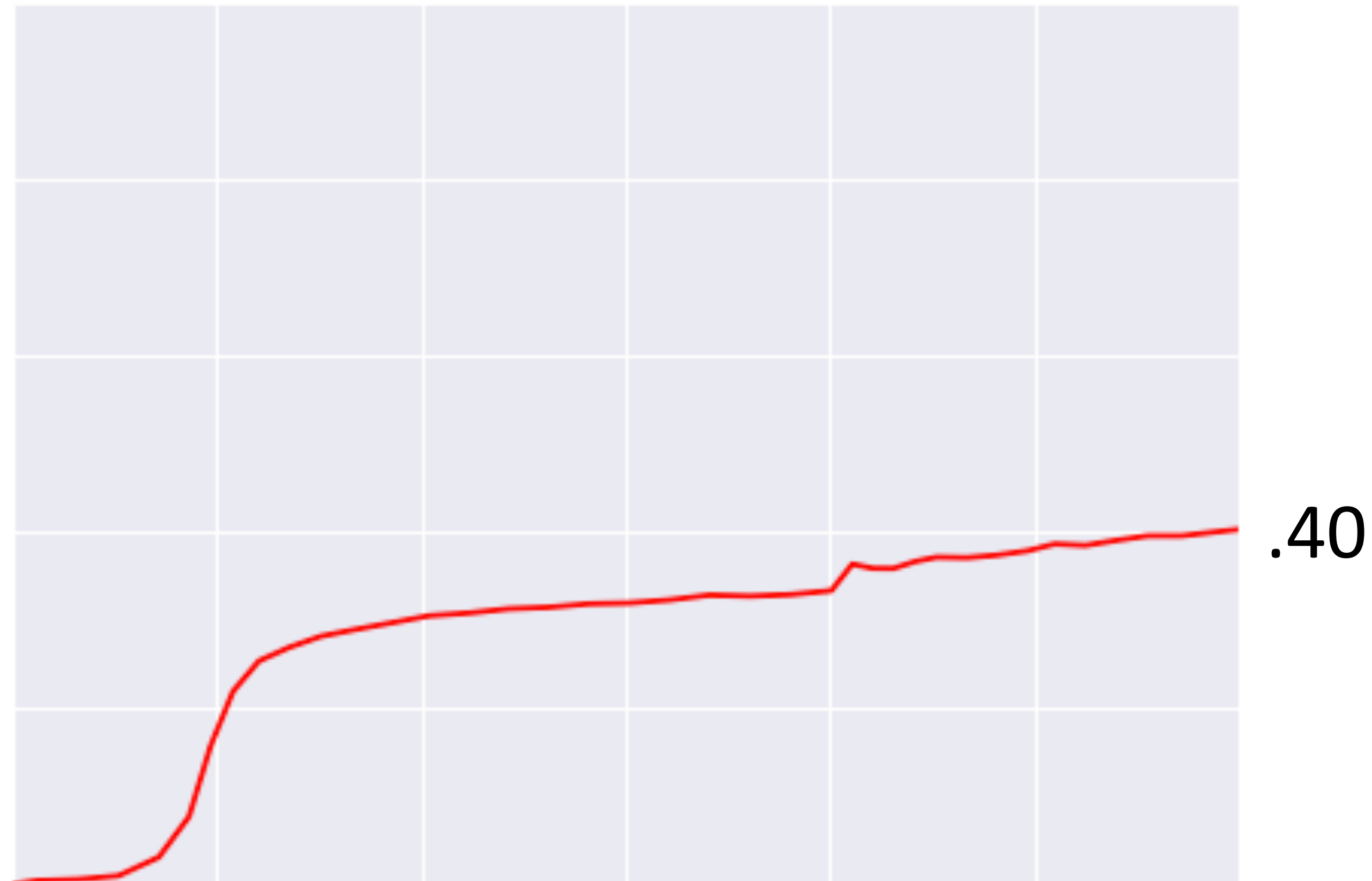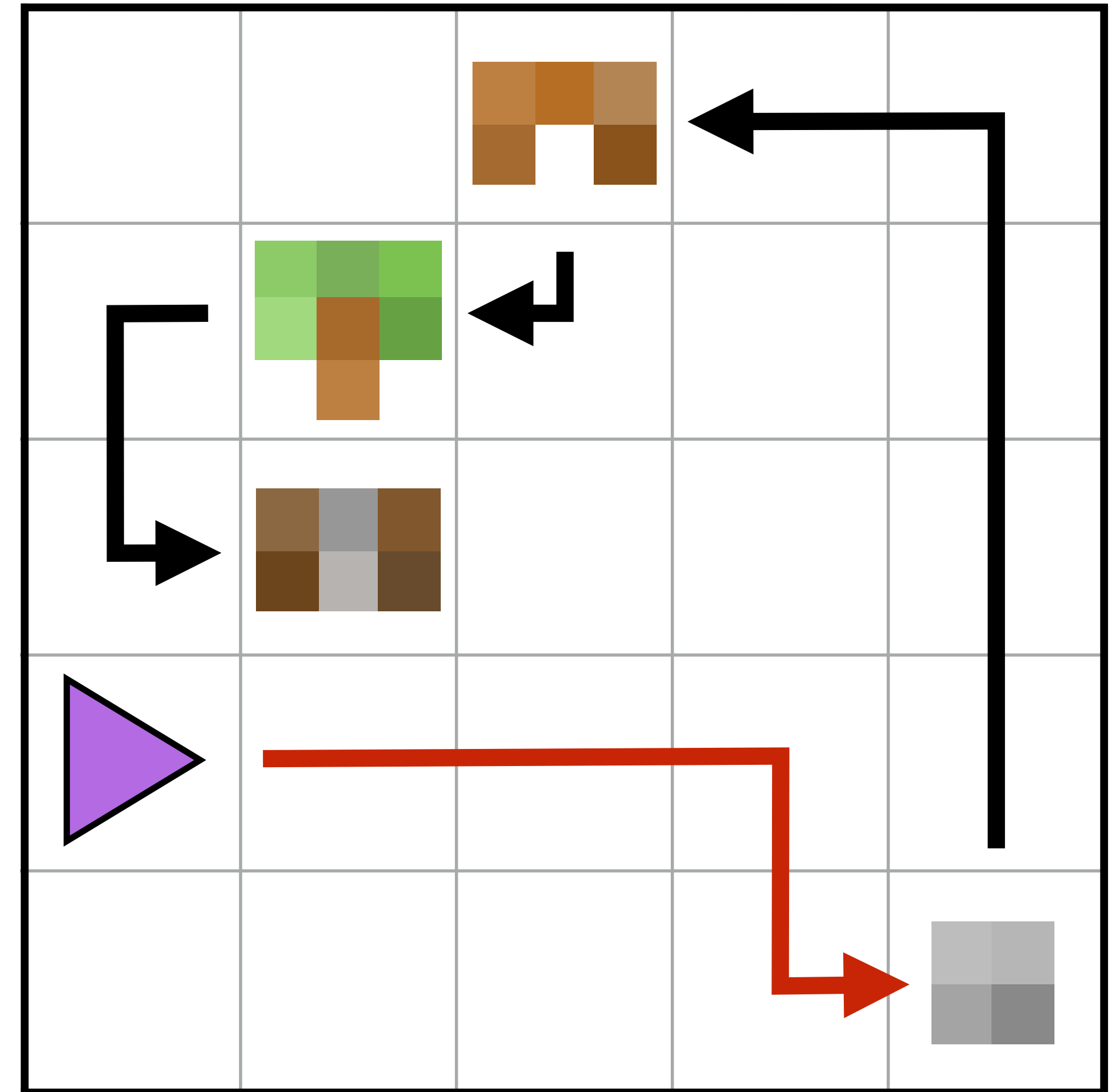
$$\left( \nabla \log \boxed{\texttt{get wood}} \right) \left( r_t - \boxed{\texttt{make nails}} \right)$$

$$\left( \nabla \log \boxed{\texttt{get iron}} \right) \left( r_t - \boxed{\texttt{make planks}} \right)$$

$$\left( \nabla \log \boxed{\texttt{get iron}} \right) \left( r_t - \boxed{\texttt{make nails}} \right)$$

# Improving policy search



Reward

.89

.40

$$\sum_{\text{tasks}} \sum_{\text{steps}} (\nabla \log \boxed{\text{SUBPOLICY}}) (r_t - \boxed{\text{TASK}})$$

# Do sketches help?

Reward

x $10^6$ episodes

Sketches: modular

Unsupervised

Sketches: joint

# The mini-craft task



Reward

x $10^6$ episodes

Sketches: modular

Sketches: joint
Unsupervised

# The cliff-walking task

# The cliff-walking task



log Reward

Sketches: modular

Sketches: joint

Unsupervised

0          1          2          3

x $10^8$ timesteps

# Zero-shot generalization

What if I see a sketch I've never seen before?

get iron

use axe

# Zero-shot generalization

What if I see a sketch I've never seen before?



Joint

**Modular**

| | Multitask | Zero-shot |
|---|---|---|
| Joint | 49 | |
| Modular | 89 | |

# Zero-shot generalization

What if I see a sketch I've never seen before?

**Joint**

**Modular**

89

77

49

1

Multitask

Zero-shot

# Fast adaptation

What if I don't get a sketch at test time?

???

# Fast adaptation

What if I don't get a sketch at test time?

**Unsupervised**

**Sketches**

47

89

Multitask                    Adaptation

# Fast adaptation

What if I don't get a sketch at test time?



Unsupervised
Sketches

| | Multitask | Adaptation |
|---|---|---|
| Unsupervised | 47 | 42 |
| Sketches | 89 | 76 |

# Conclusions

# A tiny bit of data goes a long way

```
make plank     get wood    use toolshed
make stick     get wood    use workbench
make cloth     get grass   use factory
make rope      get grass   use toolshed
make bridge    get iron    get wood       use factory
make bed*      get wood    use toolshed   get grass    use workbench
make axe*      get wood    use workbench  get iron     use toolshed
make shears    get wood    use workbench  get iron     use workbench
get gold       get iron    get wood       use factory  use bridge
get gem        get wood    use workbench  get iron     use toolshed   use axe
```
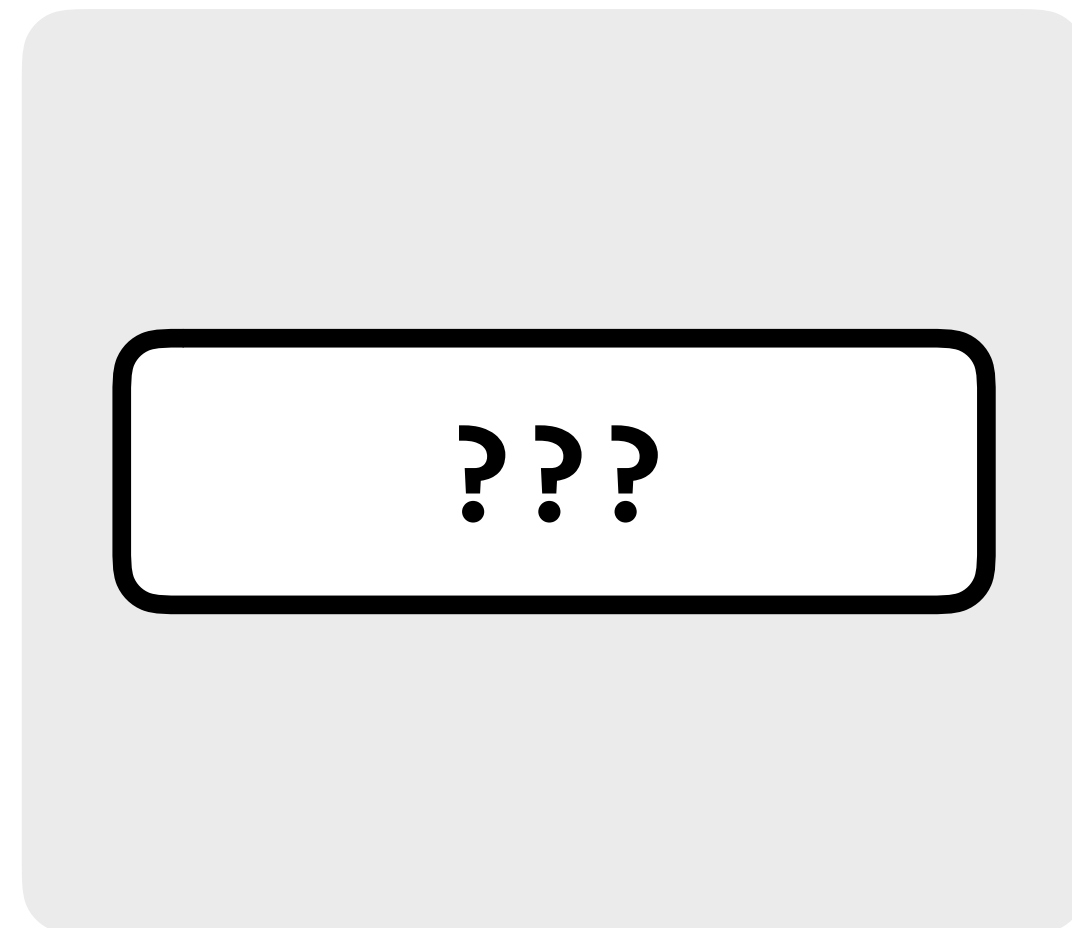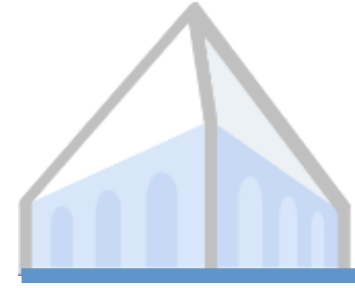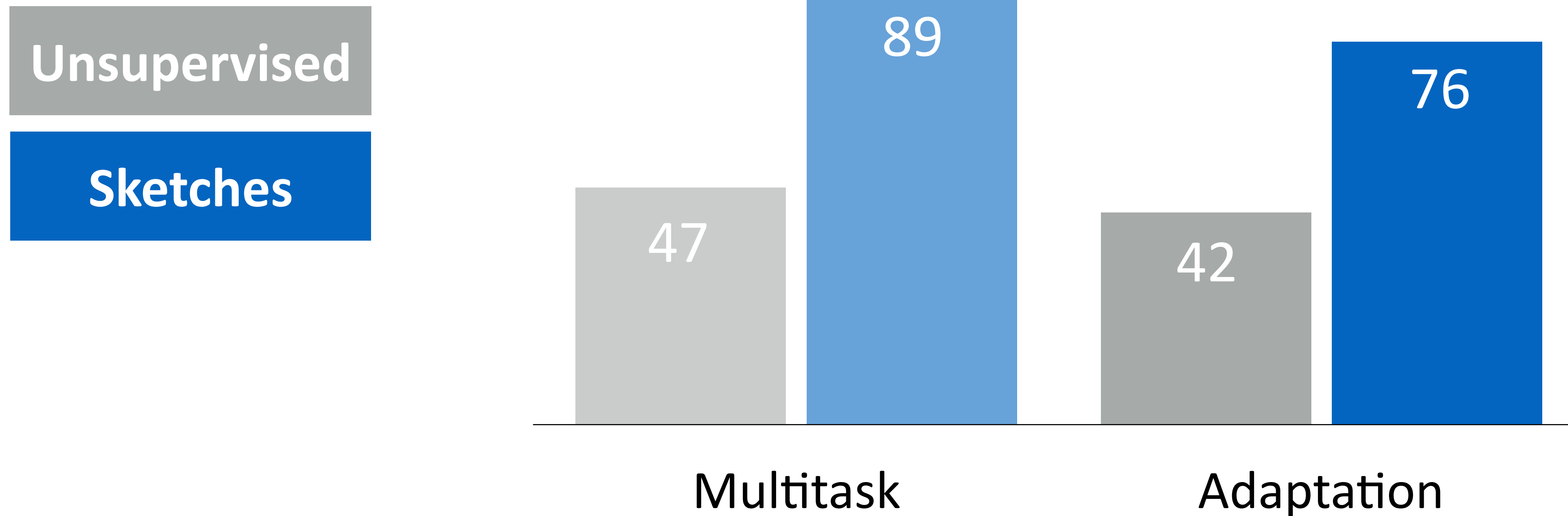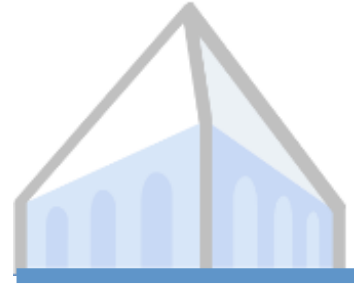
```
make plank    get
make stick    get
make cloth    get
make rope     get
make bridge   get
make bed*     get          rkbench
make axe*     get          olshed
make shears   get          rkbench
get gold      get          idge
get gem       get          olshed   use axe
```

# Thank you!

https://github.com/jacobandreas/psketch