

Quantitative Reasoning and a Bayesian View of Synthesis

with inputs on slides from Armando Solar-Lezama

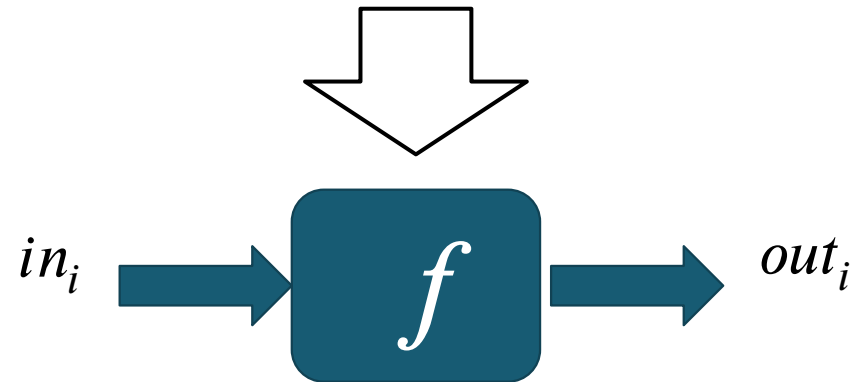
Bayes Theorem

$$P(A | B) = \frac{P(B | A) P(A)}{P(B)}$$

$$P(A \cap B) = P(B | A) P(A) = P(A | B) P(B)$$

Programming by Example

$[(in_0, out_0), (in_1, out_1), \dots, (in_k, out_k)]$



Problem is hopelessly underspecified

- Many semantically distinct programs can satisfy the examples

Bayesian View of PBE

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

$$P(f | evidence) = \frac{P(evidence | f)P(f)}{P(evidence)}$$

I/O Examples

Background
Prob/Prior

For our purpose of finding an optimal f , we can ignore $P(evidence)$ in the denominator

$P(evidence) \neq 0$

Bayesian View of PBE

Find the best f given the evidence

$$P(f \mid evidence) = \frac{P(evidence \mid f) P(f)}{P(evidence)}$$

Find the best f given the evidence

$$\begin{aligned} \operatorname{argmax}_f P(f \mid evidence) &= \operatorname{argmax}_f \frac{P(evidence \mid f) P(f)}{P(evidence)} \\ &= \operatorname{argmax}_f P(evidence \mid f) P(f) \end{aligned}$$

WARNING: $P(evidence)$ better not be zero!

$P(\textit{evidence} \mid f)$

z is a normalization constant, crucial for making sure these are probabilities, but unimportant from the point of view of optimization.

Perfectly captured I/O examples

- $P(\textit{evidence} \mid f) = P\left(\left[(in_i, out_i)\right]_i \mid f\right) = \begin{cases} 1/z & \forall_i f(in_i) = out_i \\ 0 & \textit{otherwise} \end{cases}$
- With a uniform $P(f)$ this reduces to finding any function that works

$P(f)$

So far we have been using a uniform P

- $$P(f) = \begin{cases} 1/Z & \text{if } f \text{ belongs to the space of programs} \\ 0 & \text{otherwise} \end{cases}$$

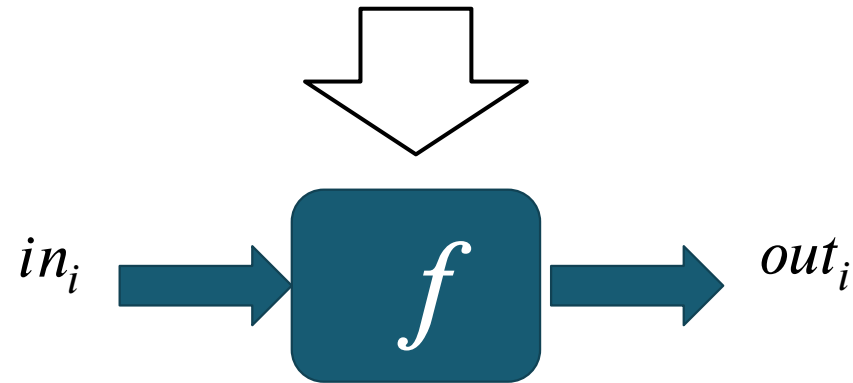
Shortest programs are better than longer programs

- $$P(f) = \begin{cases} \frac{1}{Z} * e^{-len(f)} & \text{if } f \text{ belongs to the space of programs} \\ 0 & \text{otherwise} \end{cases}$$

Could we learn $P(f)$?

Programming by Example

$[(in_0, out_0), (in_1, out_1), \dots, (in_k, out_k)]$



Problem is hopelessly underspecified

- Many semantically distinct programs can satisfy the examples

$$P(f \mid [(in_i, out_i)]_i) \approx P_f(f) * P_{io} \left([(in_i, out_i)]_i \mid f \right)$$

$P(\textit{evidence} \mid f)$ for Synthesis under errors

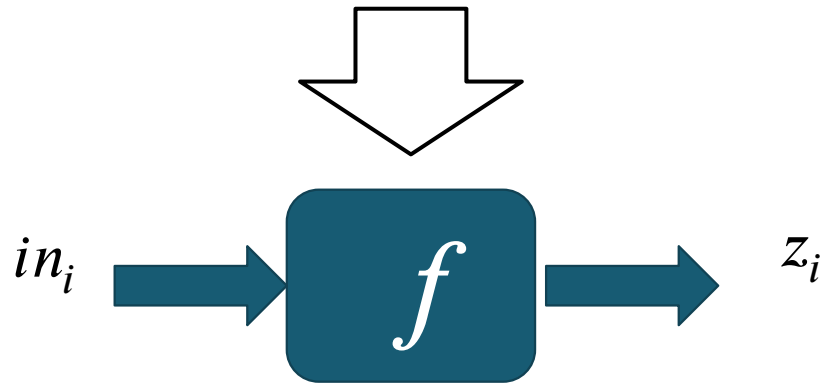
$$P(A \mid B) = \frac{P(B \mid A)P(A)}{P(B)}$$

Imperfectly captured independent I/O examples

- $P(\textit{evidence} \mid f) = P\left(\left[(in_i, out_i)\right]_i \mid f\right) = \prod_i P_{o|z}(out_i \mid f, in_i)P(in_i)$
- For the purposes of maximizing $P(f)$, $P(in_i)$ can be ignored if all inputs are equally likely

Learning from noisy data

$[(in_0, out_0), (in_1, out_1), \dots, (in_k, out_k)]$



Need to trade off quality of f against faithfulness to data

- This requires an error model

$$P(f \mid [(in_i, out_i)]_i) \approx P_f(f) * \prod_i P_{o|z}(out_i \mid f, in_i)$$

Example: Off-by-one Errors

Suppose we know off-by-one errors are possible in the data

$$\bullet P_{o|z}(out_i \mid f, in_i) = \begin{cases} 0.5 & f(in_i) = out_i \\ 0.25 & f(in_i) = out_i \pm 1 \\ 0 & else \end{cases}$$

If $p(f)$ is uniform, this reduces to

- “Discard programs that are more than one-off on any input”
- “From the remaining programs, select the one that matches the most examples”

Off-by-one Errors

Suppose we know off-by-one errors are possible in the data but others are possible as well.

$$\bullet P_{o|z}(out_i \mid f, in_i) = \begin{cases} \frac{1}{Z} 0.5 & f(in_i) = out_i \\ \frac{1}{z} 0.25 & f(in_i) = out_i \pm 1 \\ \epsilon & \textit{else} \end{cases}$$

Non-uniform $P(f)$

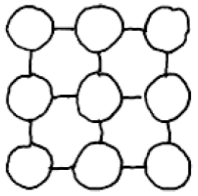
Trade off $P_{o|z}(out_i \mid f, in_i)$ against $P(f)$

- A solution that misses more outputs may still be preferable if it has much higher probability

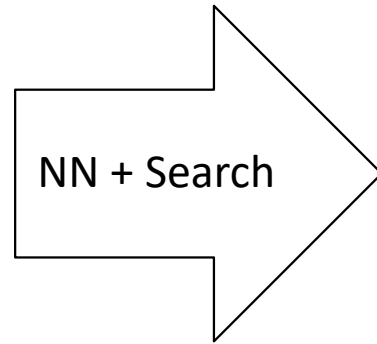
Learning to Infer Graphics Programs from Hand-Drawn Images

Kevin Ellis, Daniel Ritchie, Josh Tenenbaum

From images to programs



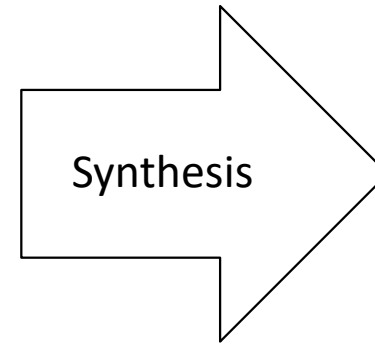
Hand Drawn Figure



```
Circle(5,8)
Circle(2,8)
Circle(8,11)
Line(2,9, 2,10)
Circle(8,8)
Line(3,8, 4,8)
Line(3,11, 4,11)
Line(8,9, 8,10)
Circle(5,14)
```

... etc. ...; 21 lines

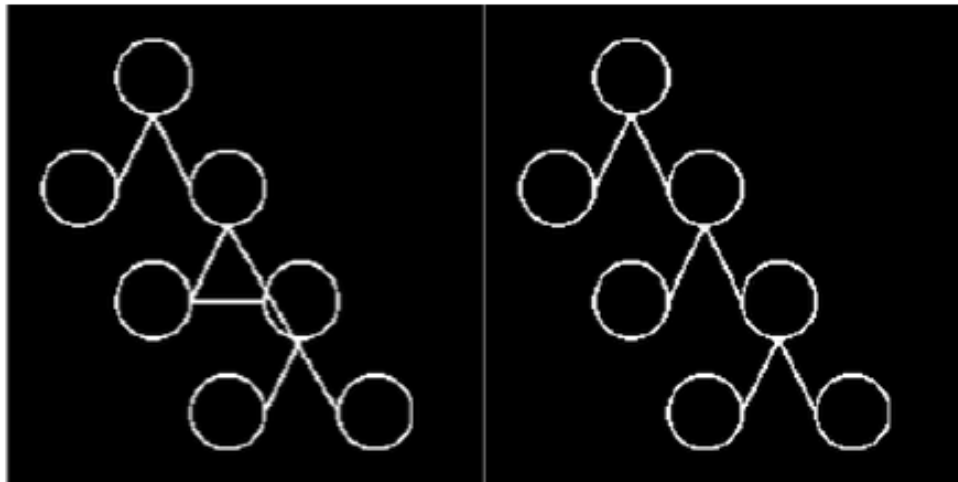
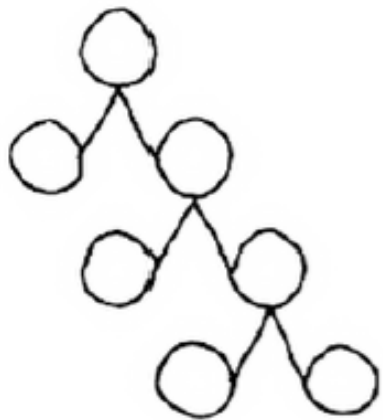
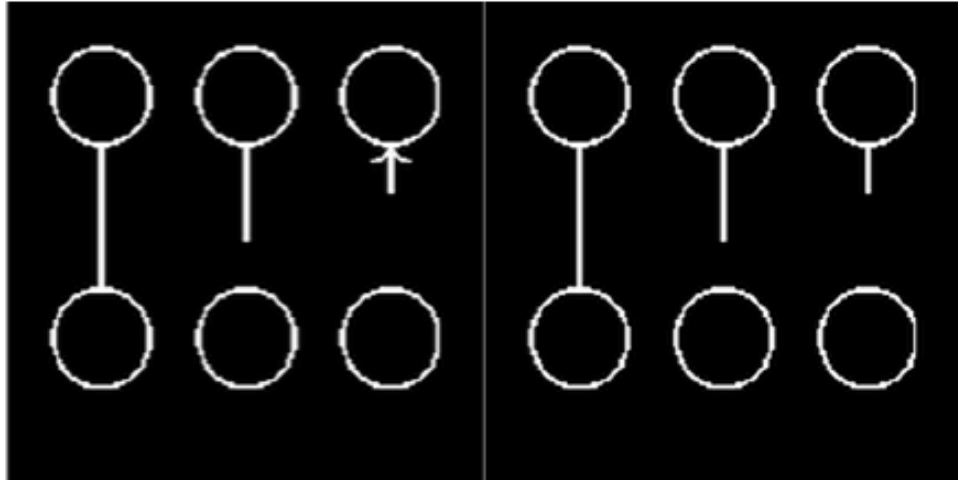
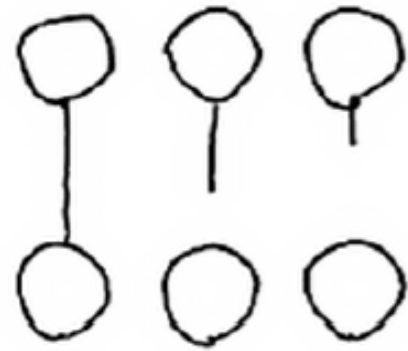
Description of
elements in the drawing



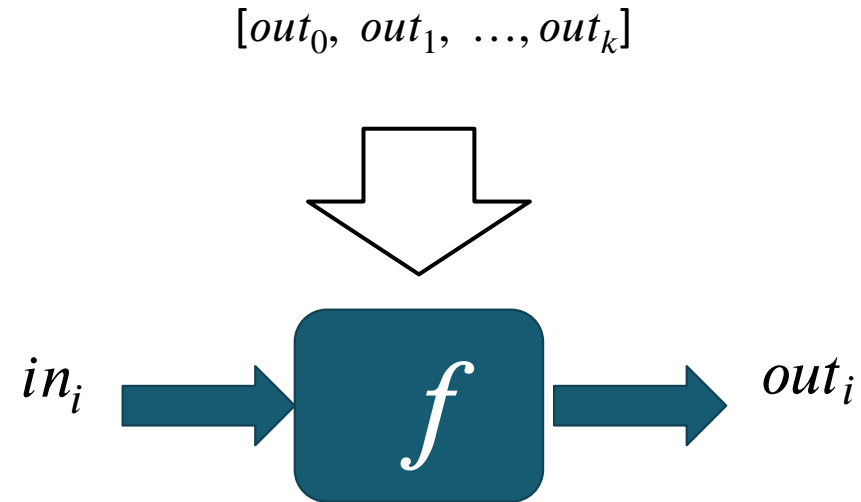
```
for(i<3)
  for(j<3)
    if(j>0)
      line(-3*j+8,-3*i+7,
           -3*j+9,-3*i+7)
    line(-3*i+7,-3*j+8,
         -3*i+7,-3*j+9)
  circle(-3*j+7,-3*i+7)
```

Program representation
of drawing

Why? Correcting errors in perception



Unsupervised learning



This is hopelessly underspecified

- Can we identify the process that generated the sequence?

Unsupervised learning

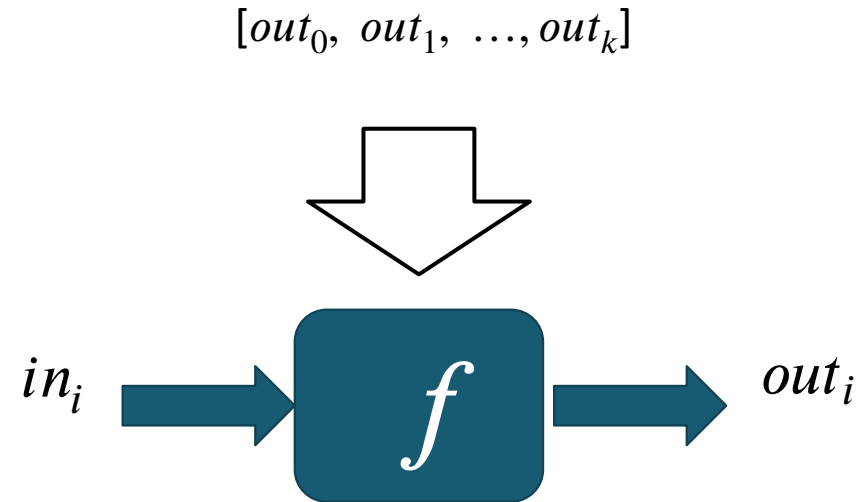
$$P\left(f, [in_i] \mid [out_i]\right) = \frac{P\left([out_i] \mid f, [in_i]\right) P(f, [in_i])}{P([out_i])}$$

Assuming independence:

$$P([out_i] \mid f, [in_i]) = \prod_i P(out_i \mid f, in_i)$$

$$P(f, [in_i]) = P(f) * \prod_i P(in_i)$$

Unsupervised learning



This is hopelessly underspecified

- Can we identify the process that generated the sequence?

$$P\left(f, [in_i]_i \mid [out_i]_i\right) \approx P_f(f) * \prod_i P_{o|z}(out_i \mid f, in_i) * P_{in}(in_i)$$

To marginalize or not to Marginalize

$$P\left(f, [in_i]_i \mid [out_i]_i\right)$$

Probability that a given function and inputs were the cause for an observed series of outputs

$$\sum_{[in_i]_i} P\left(f, [in_i]_i \mid [out_i]_i\right) P([in_i]_i)$$

Probability that a given function is consistent with the observed outputs

Which of the two functions above should you be optimizing?

- Formulation on the left is easier to solve for
 - especially with symbolic methods

Maximum Likelihood vs Sampling

Often your goal is to find the most likely f

- $\max_f P_f(f \mid \dots)$

For some situations, sampling from $P_f(f \mid \dots)$ is better

- The most likely is not necessarily the one you want
- E.g. in PBE the function the user has in mind may not be the “best”

Inversion, Rejection, Relationship and Approximation

Isn't there a whole field looking into this?

Machine learning has been studying these problems for a while

What we bring to the table:

- Flexible spaces of functions
- Complex distributions
- Powerful symbolic search techniques

Machine Learning for Synthesis

- DeepCoder : Learning to Write Programs,
 - Balog et al. ICLR 2017
- DreamCoder : Bootstrapping Inductive Program Synthesis with Wake-Sleep Library Learning.
 - Ellis et al. PLDI 2021

DreamCoder

Kevin Ellis, Catherine Wong, Maxwell Nye, Mathias
Sable-Meyer, Luc Cary, Lucas Morales, Like Hewitt, Armando
Solar-Lezama, Joshua B. Tenenbaum

Whats happening here?

3 1 2 5 → 1 2 3 5

9 4 7 2 3 → 2 3 4 7 9

8 5 1 → 1 5 8

6 4 2 5 → ?

Sorting: This core feature of Humans is hard for ML

Combine ML techniques with Abstraction Learning

Program Synthesis (Inductive Program Synthesis)

Goal: acquire domain-specific knowledge needed to induce a class of programs

- Library of abstractions (domain specific language) Assumed this is given
- Inference strategy (synthesis algorithm)

Library Learning

Initial Primitives

·
·
map
fold
if
cons
>
·
·

Sample Problem: sort list

[9 2 7 1] → [1 2 7 9]
[3 8 9 4 2] → [2 3 4 8 9]
[6 2 2 3 8 5] → [2 2 3 5 6 8]
...

Library Learning

Initial Primitives

·
·
map
fold
if
cons
>
·
·

Sample Problem: sort list

[9 2 7 1] → [1 2 7 9]
[3 8 9 4 2] → [2 3 4 8 9]
[6 2 2 3 8 5] → [2 2 3 5 6 8]
...

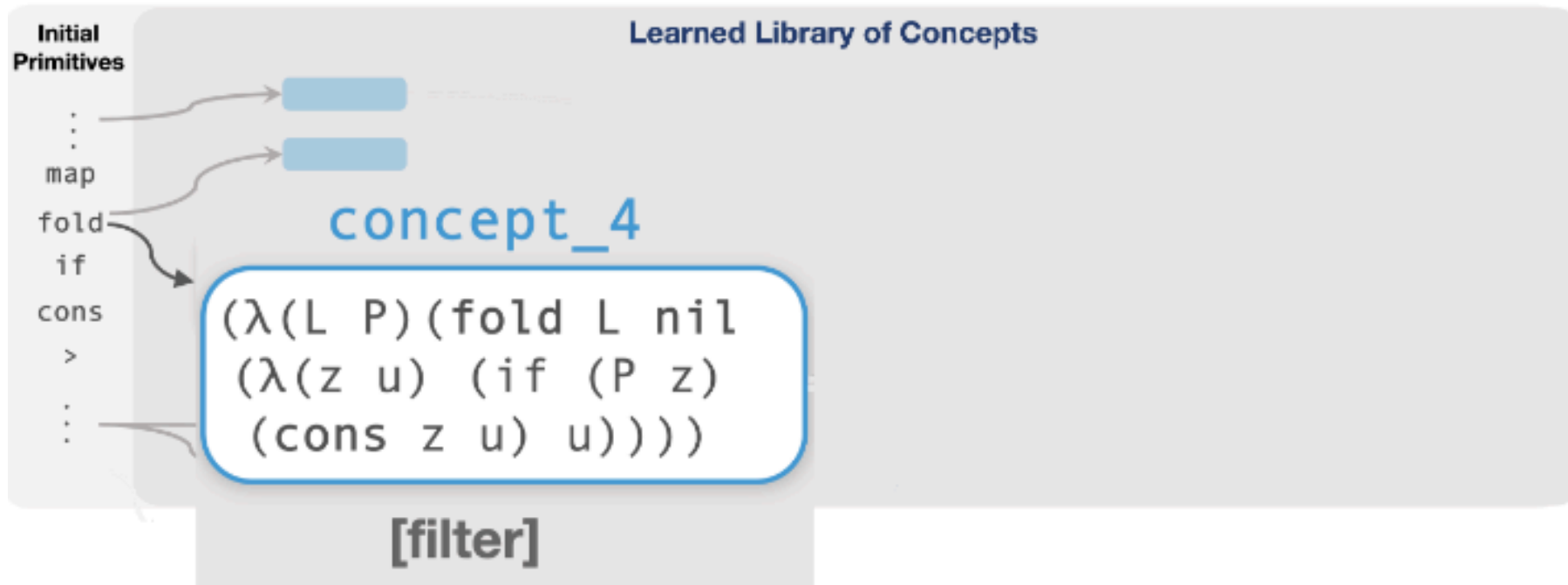
Library Learning



Sample Problem: `sort list`

```
[9 2 7 1] → [1 2 7 9]
[3 8 9 4 2] → [2 3 4 8 9]
[6 2 2 3 8 5] → [2 2 3 5 6 8]
...
```

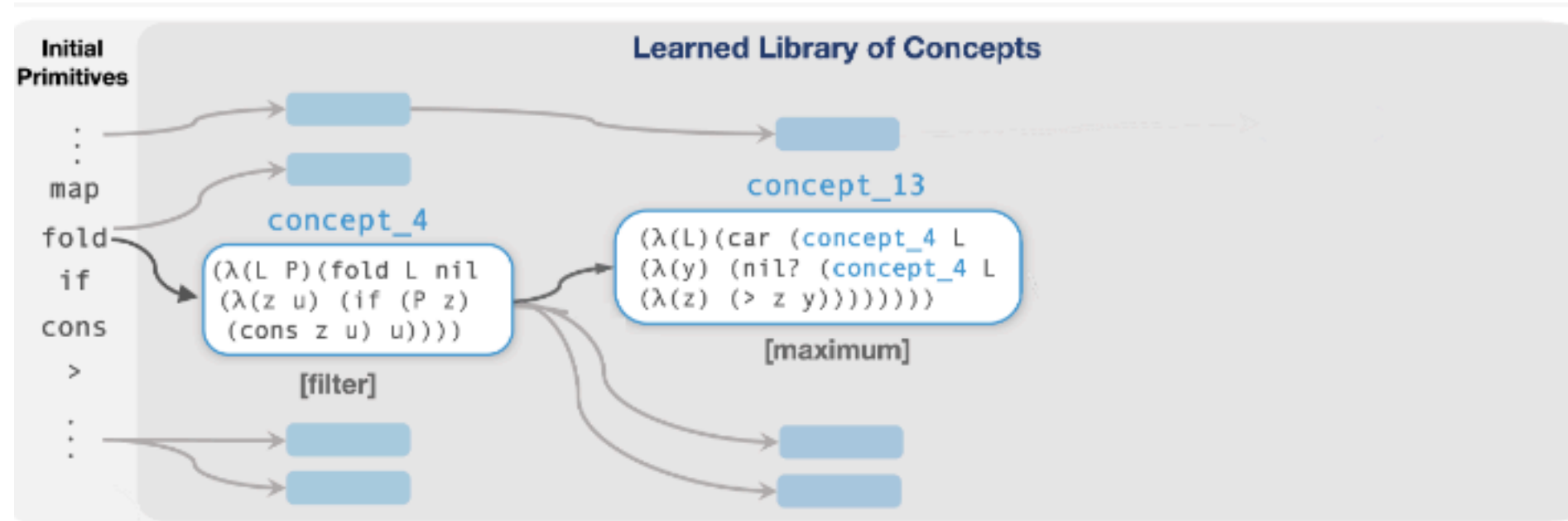
Library Learning



Sample Problem: sort list

[9 2 7 1] → [1 2 7 9]
[3 8 9 4 2] → [2 3 4 8 9]
[6 2 2 3 8 5] → [2 2 3 5 6 8]
...

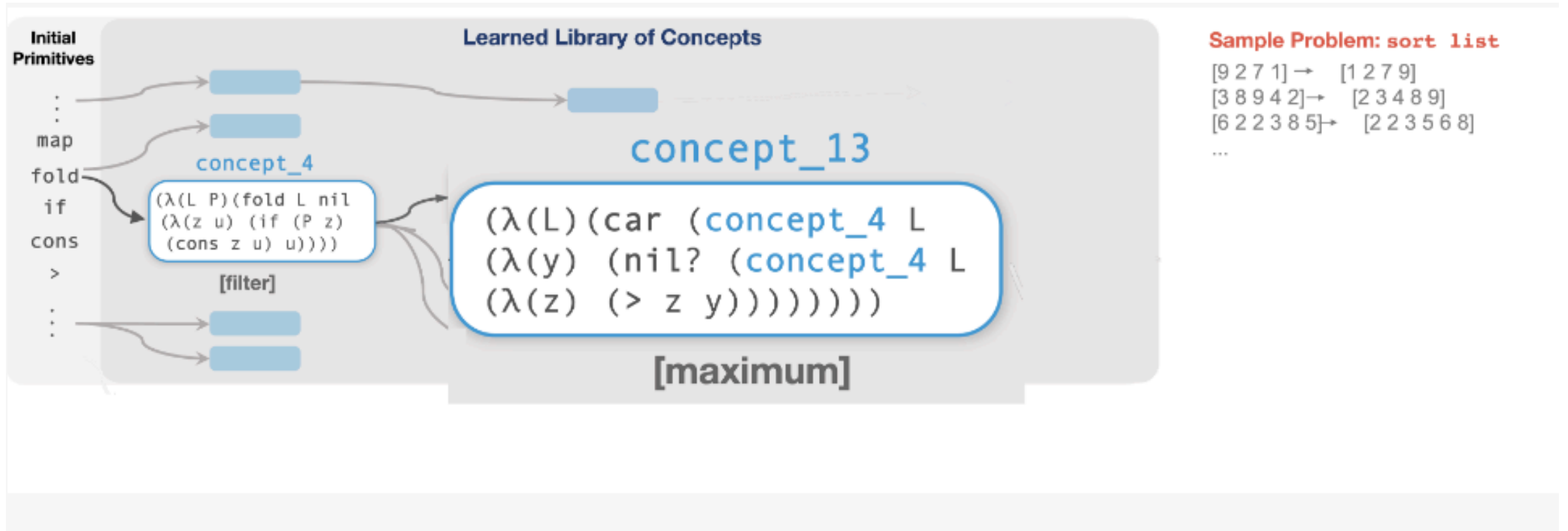
Library Learning



Sample Problem: sort list

[9 2 7 1] → [1 2 7 9]
[3 8 9 4 2] → [2 3 4 8 9]
[6 2 2 3 8 5] → [2 2 3 5 6 8]
...

Library Learning

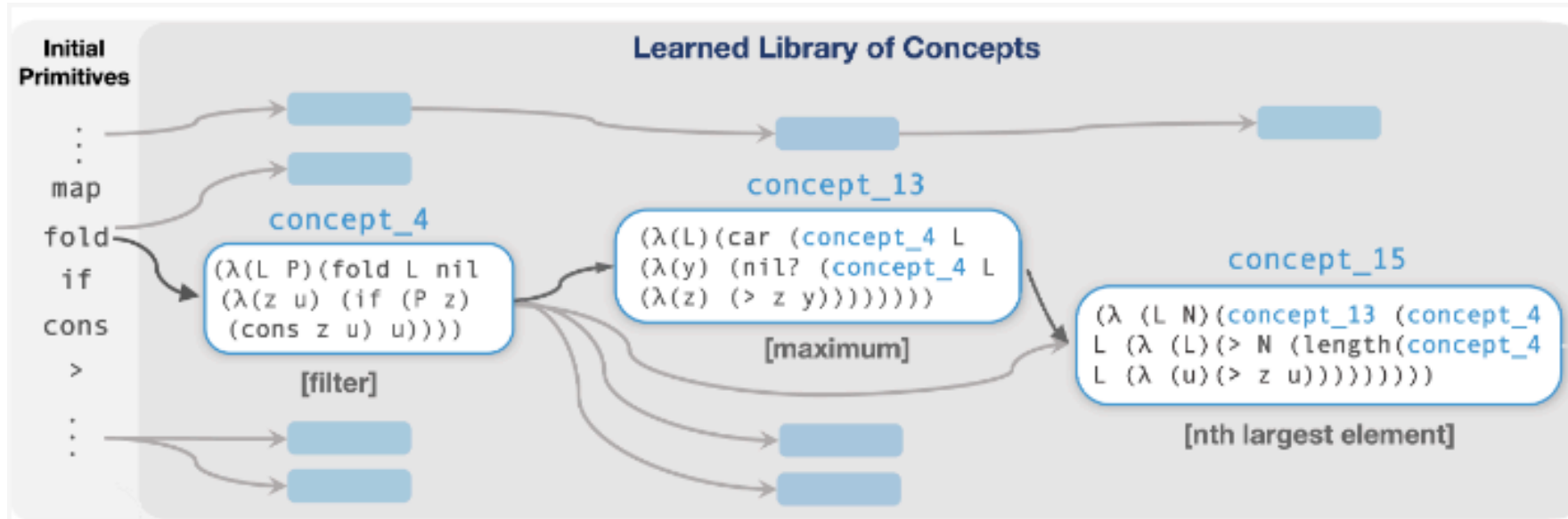


Sample Problem: sort list

```
[9 2 7 1] → [1 2 7 9]
[3 8 9 4 2] → [2 3 4 8 9]
[6 2 2 3 8 5] → [2 2 3 5 6 8]
```

...

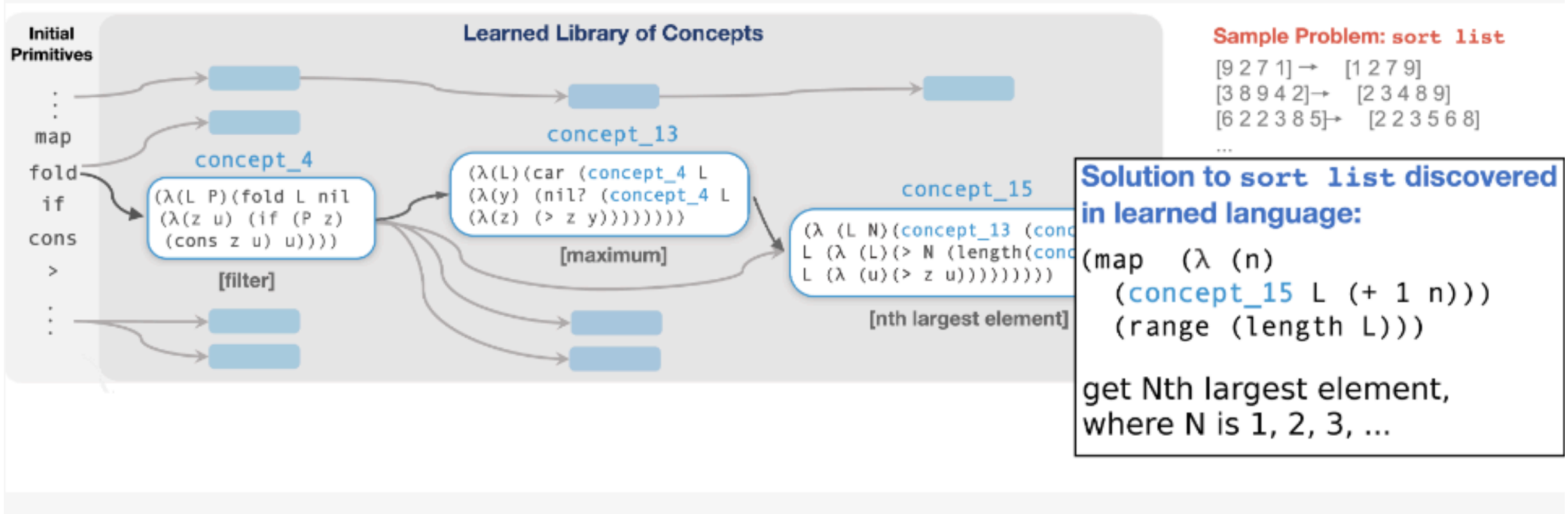
Library Learning



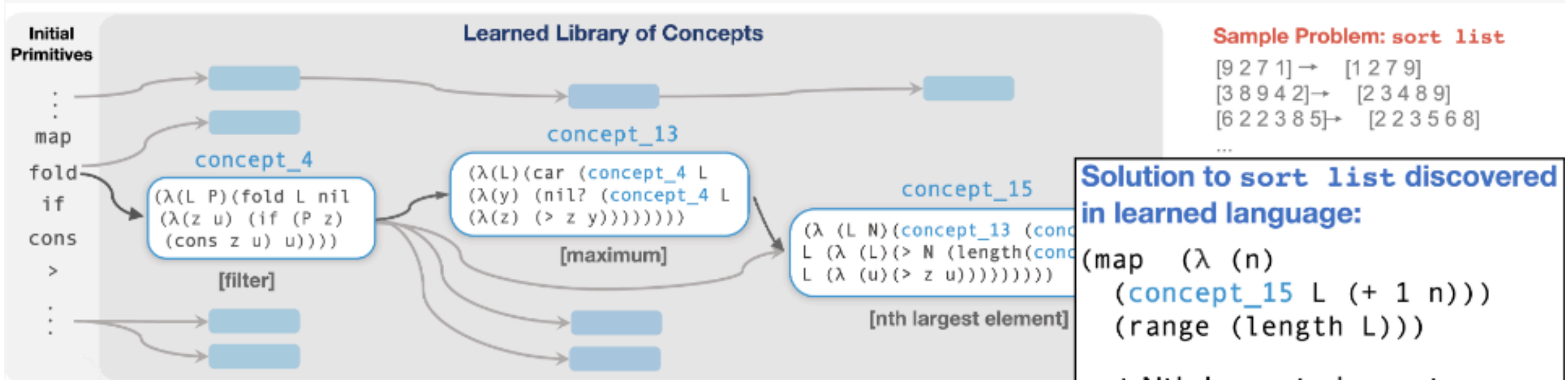
Sample Problem: sort list

```
[9 2 7 1] → [1 2 7 9]
[3 8 9 4 2] → [2 3 4 8 9]
[6 2 2 3 8 5] → [2 2 3 5 6 8]
...
```

Library Learning



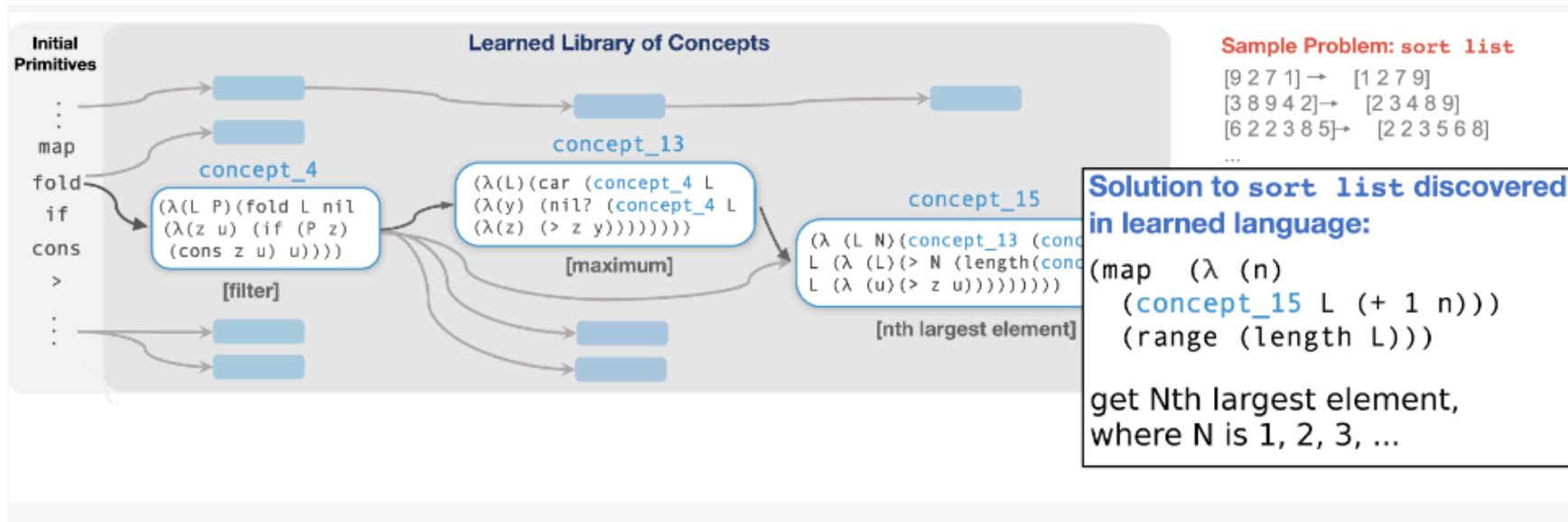
Library Learning



Solution rewritten in initial primitives:

```
(lambda (x) (map (lambda (y) (car (fold (fold x nil (lambda (z u) (if (gt? (+ y 1) (length (fold x nil (lambda (v w) (if (gt? z v) (cons v w) w)))) (cons z u) u)) nil (lambda (a b) (if (nil? (fold (fold x nil (lambda (c d) (if (gt? (+ y 1) (length (fold x nil (lambda (e f) (if (gt? c e) (cons e f) f)))) (cons c d) d))) nil (lambda (g h) (if (gt? g a) (cons g h) h)))) (cons a b) b)))) (range (length x))))
```

Library Learning



Solution rewritten in initial primitives:

```
(lambda (x) (map (lambda (y) (car (fold (fold x nil (lambda (z u) (if (gt? (+ y 1) (length (fold x nil (lambda (v w) (if (gt? z v) (cons v w) w)))) (cons z u) u))) nil (lambda (a b) (if (nil? (fold (fold x nil (lambda (c d) (if (gt? (+ y 1) (length (fold x nil (lambda (e f) (if (gt? c e) (cons e f) f)))) (cons c d) d))) nil (lambda (g h) (if (gt? g a) (cons g h) h)))) (cons a b) b)))) (range (length x))))
```

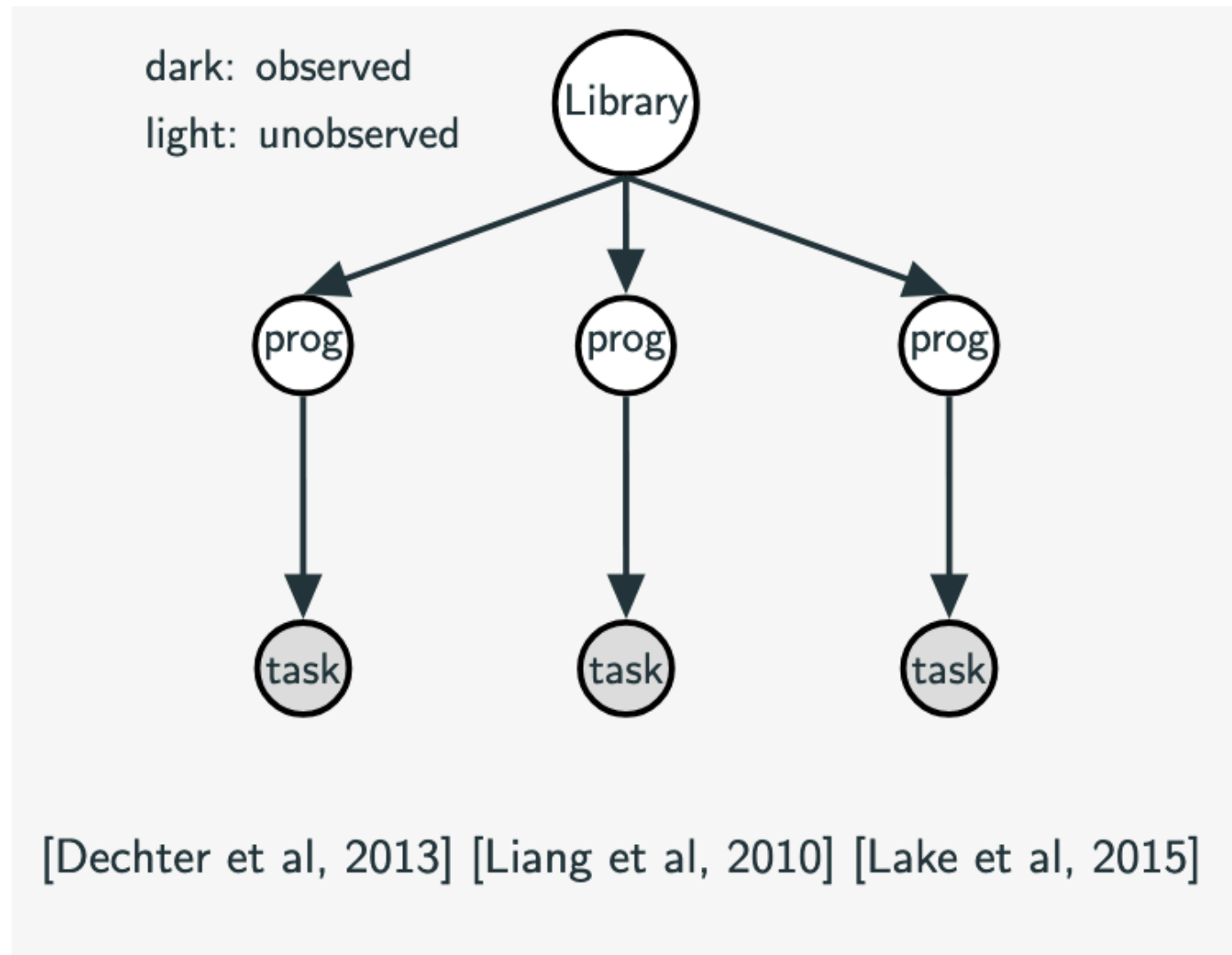
induced sort program found in ≤ 10 min. Brute-force search without learned library would take $\approx 10^{73}$ years

DreamCoder

- **Wake:** Solve problems by writing programs
- **Sleep:** Improve library and neural recognition model:
 - **Abstraction sleep:** Improve library
 - **Dream sleep:** Improve neural recognition model

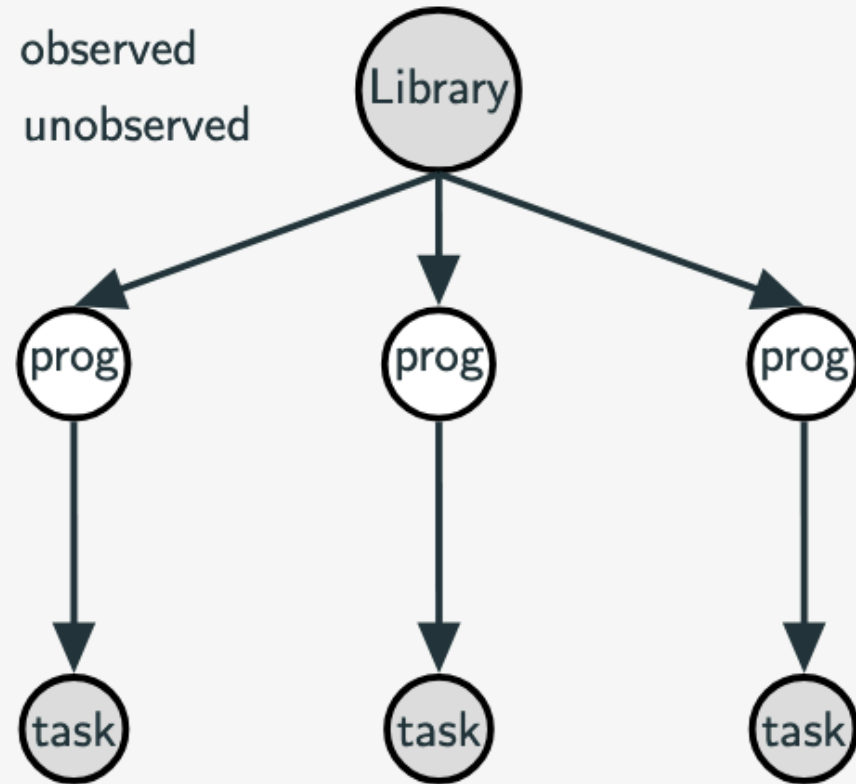
cf. Helmholtz machine, wake/sleep neural network training algorithms

Library Learning as Bayesian inferences



Library Learning as Bayesian inferences

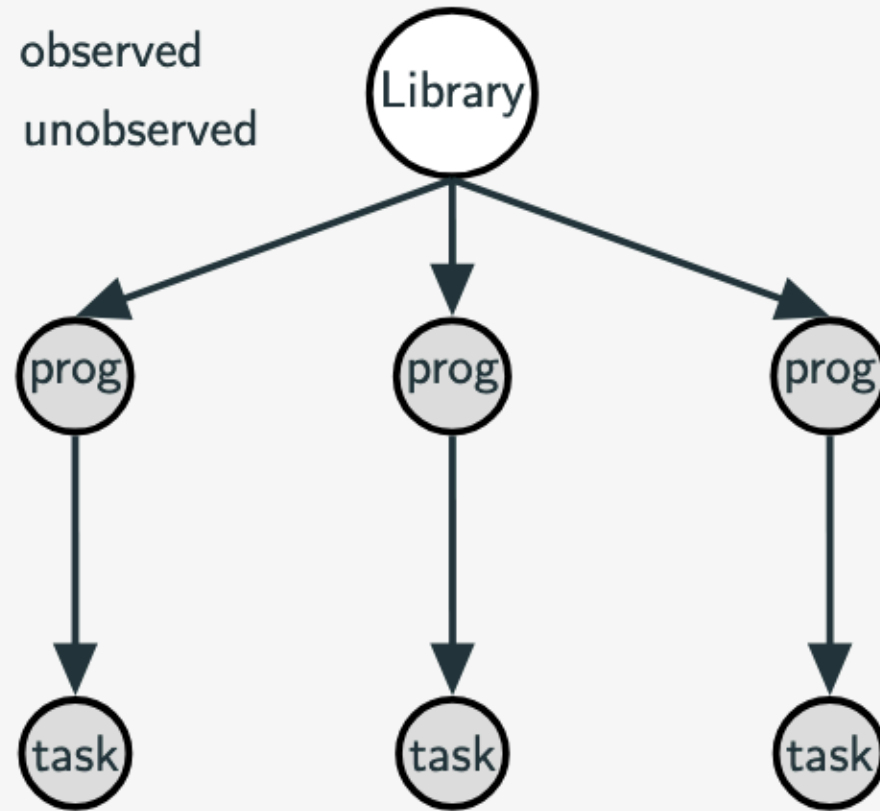
dark: observed
light: unobserved



[Dechter et al, 2013] [Liang et al, 2010] [Lake et al, 2015]

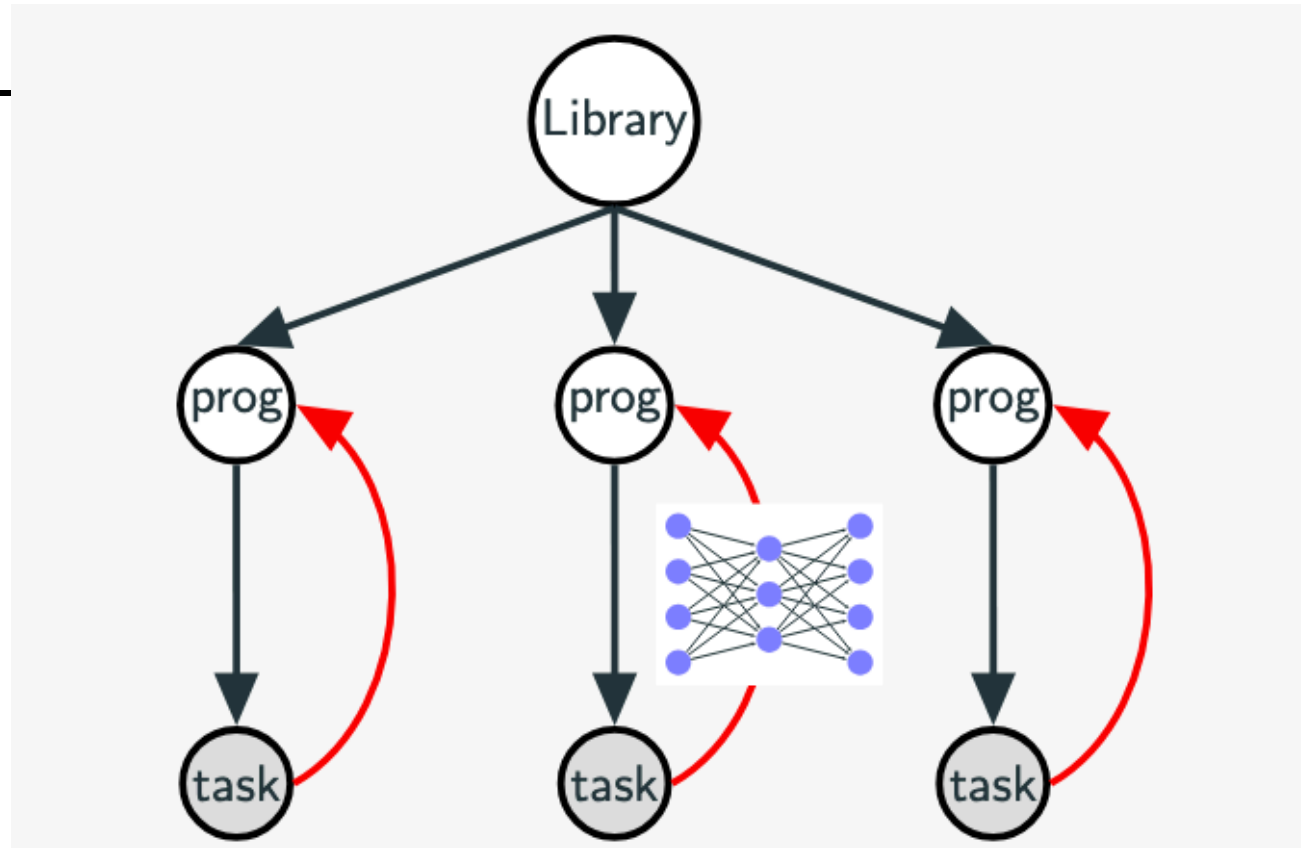
Library Learning as Bayesian inferences

dark: observed
light: unobserved

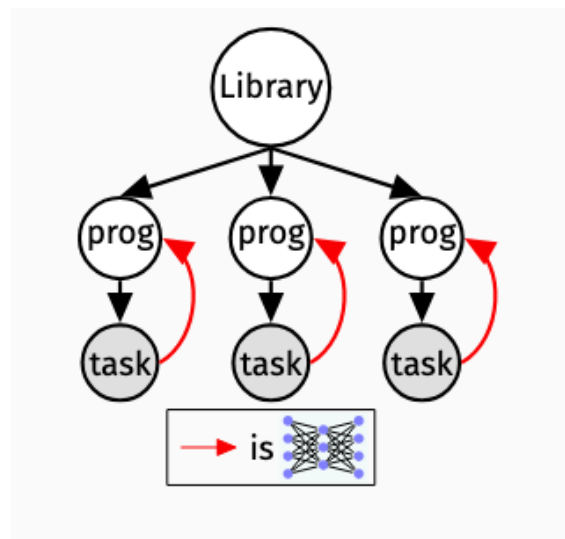


[Dechter et al, 2013] [Liang et al, 2010] [Lake et al, 2015]

Library Learning as neurally-guided Bayesian inference



library learning via program analysis +
new neural inference network for program synthesis +
better program representation (Lisp+polymorphic types [Milner 1978])



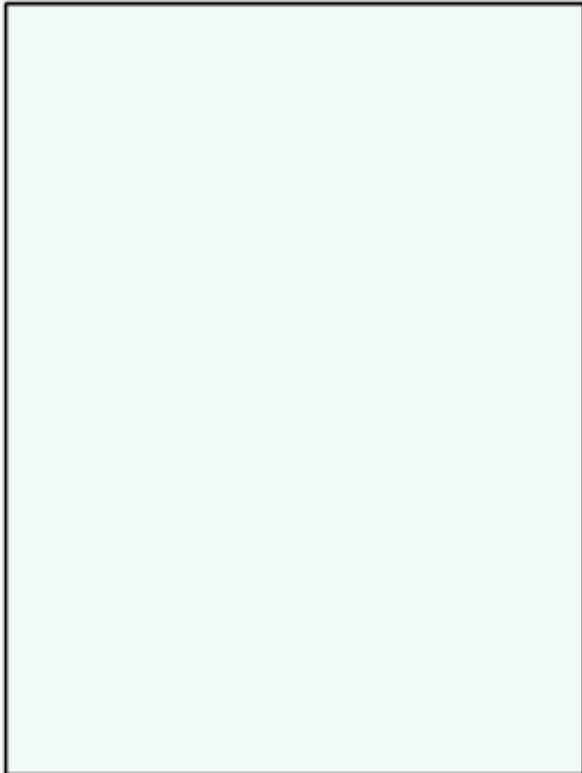
WAKE



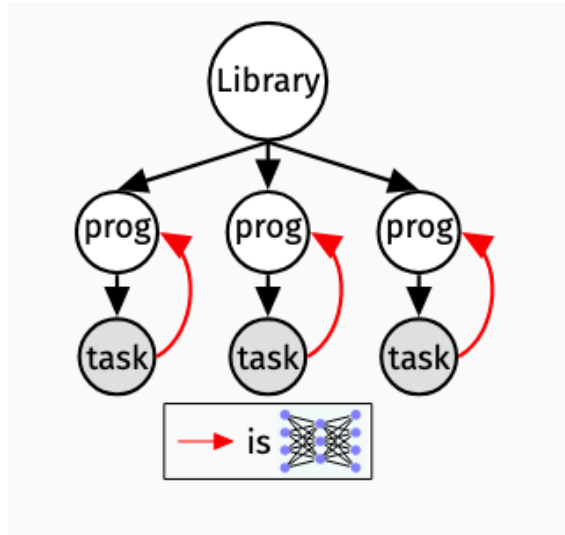
Synthesis

Train the recognition model

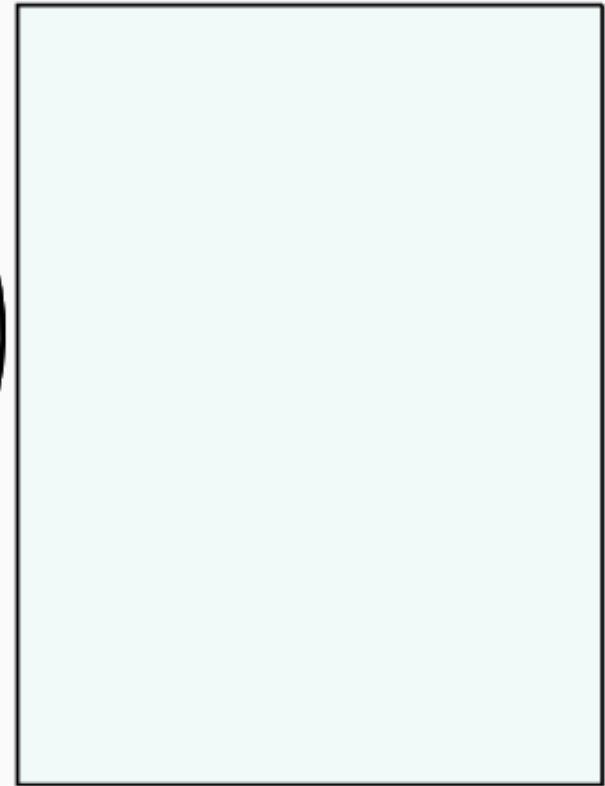
SLEEP: ABSTRACTION



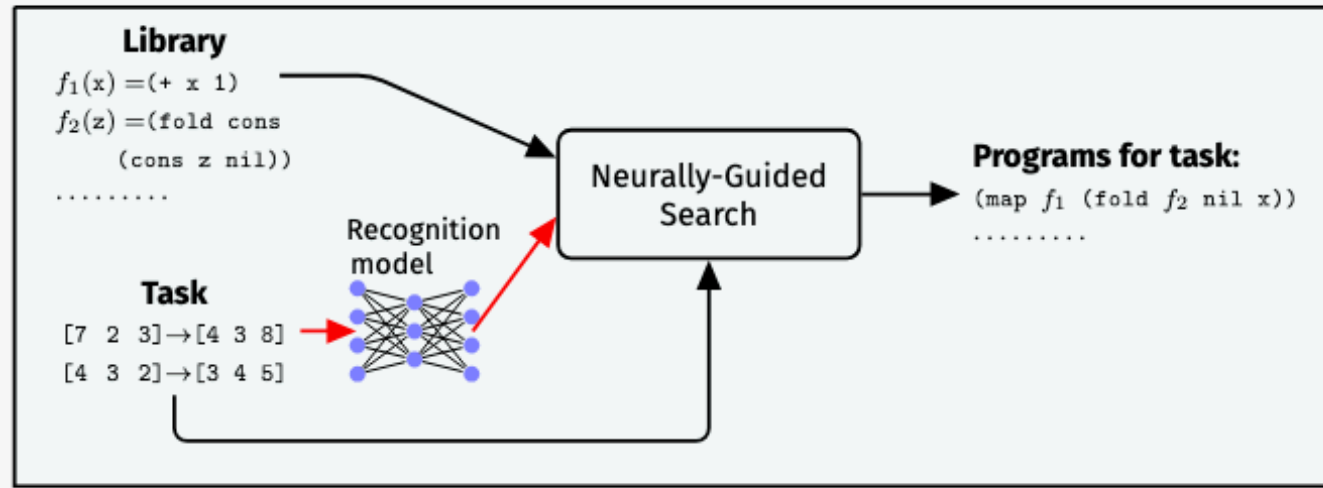
Identify re-used programs



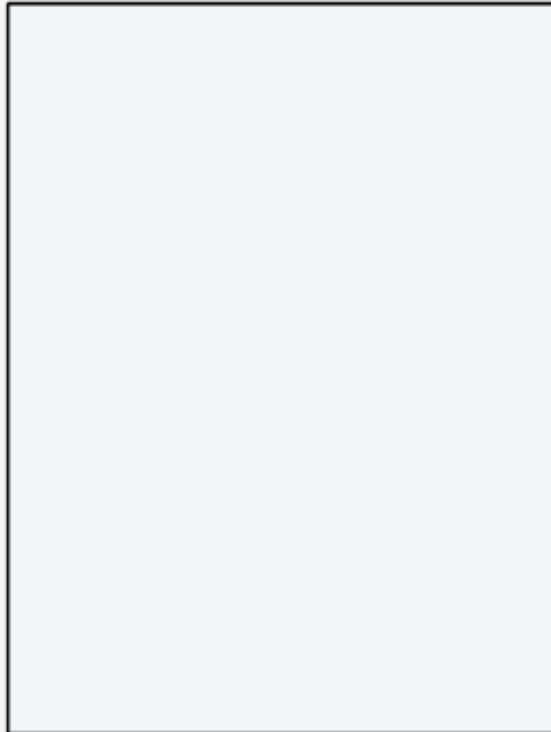
SLEEP: DREAMING



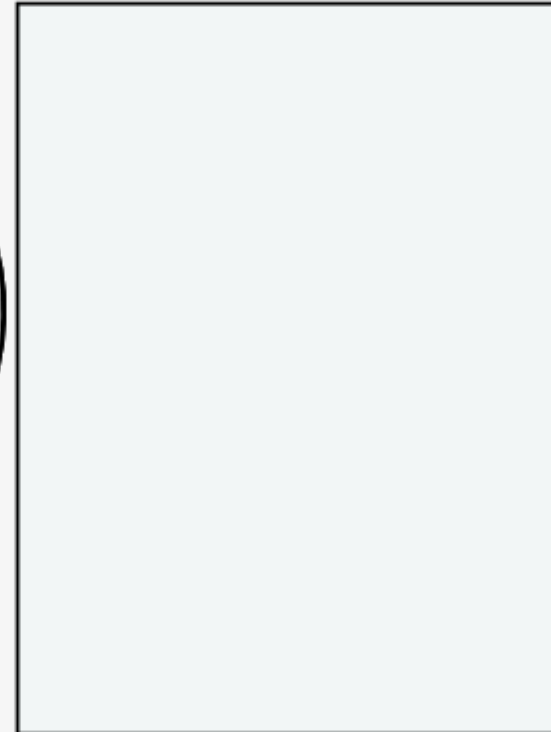
WAKE



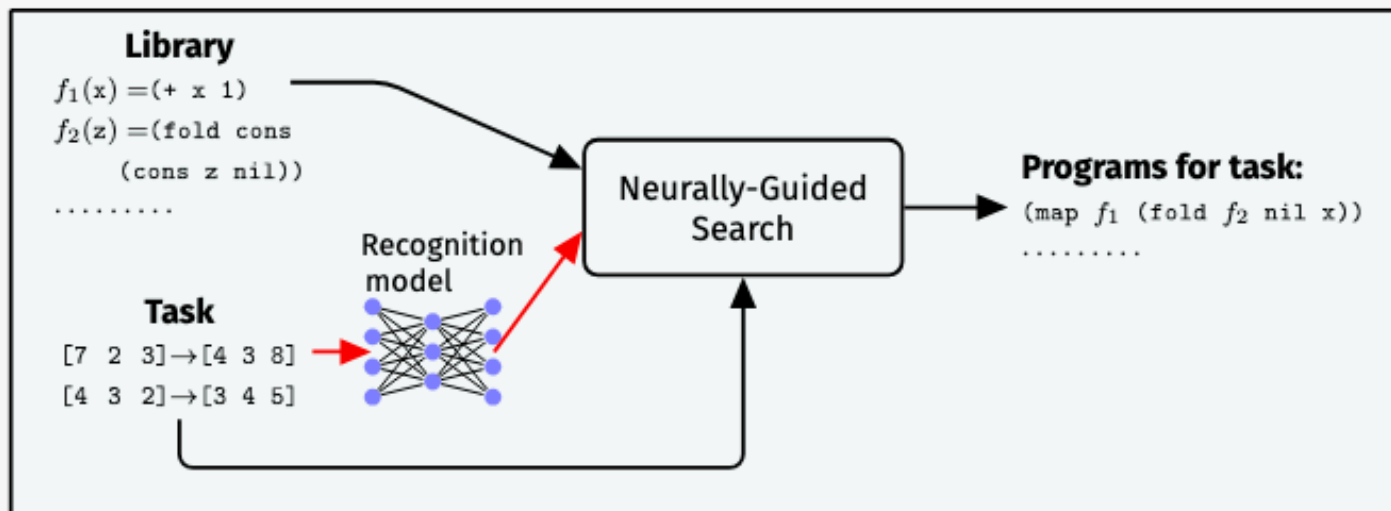
SLEEP: ABSTRACTION



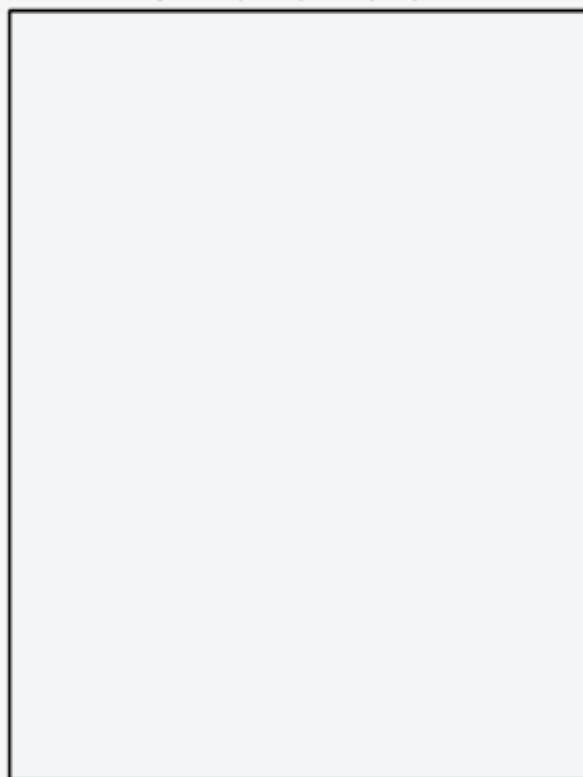
SLEEP: DREAMING



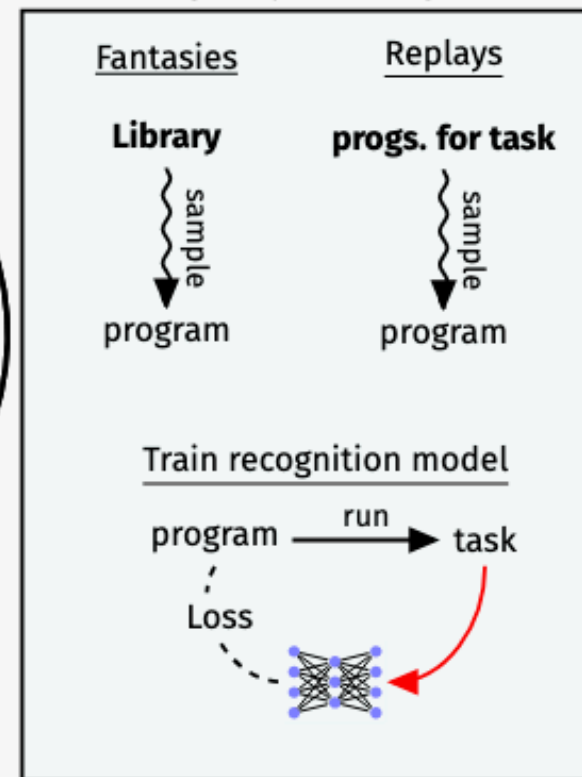
WAKE



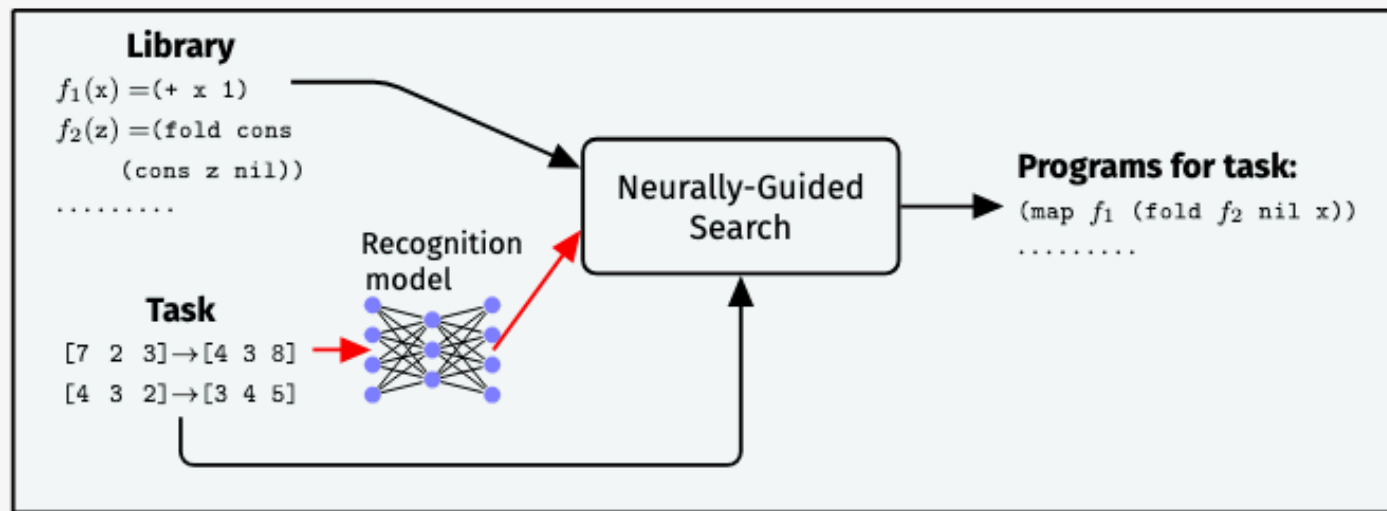
SLEEP: ABSTRACTION



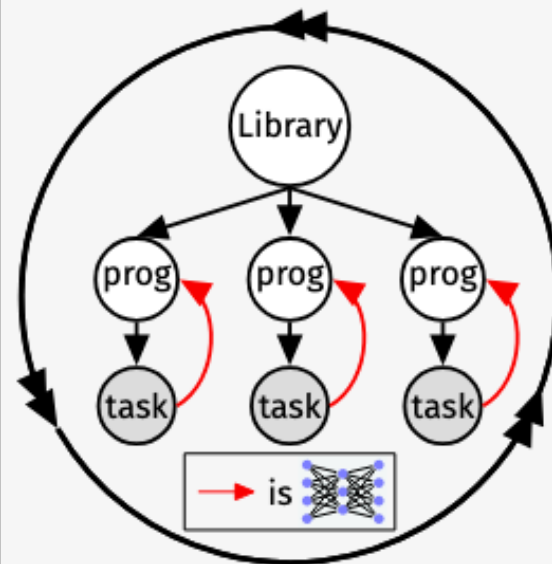
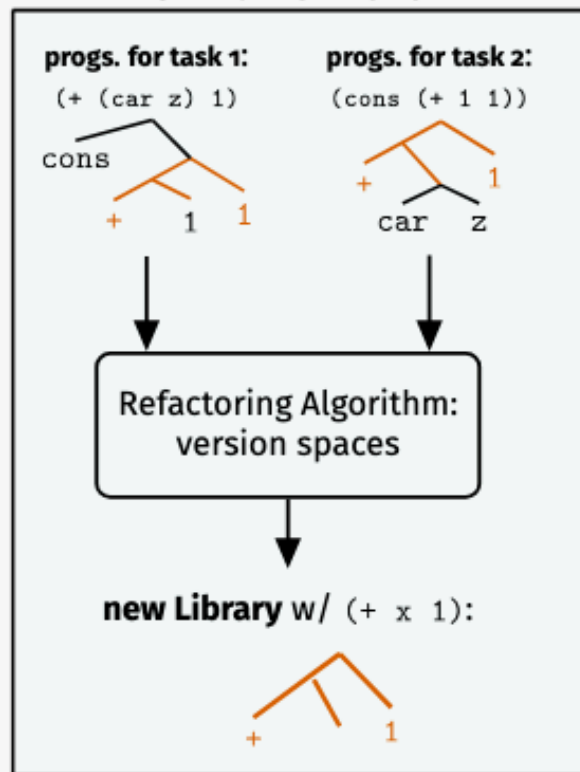
SLEEP: DREAMING



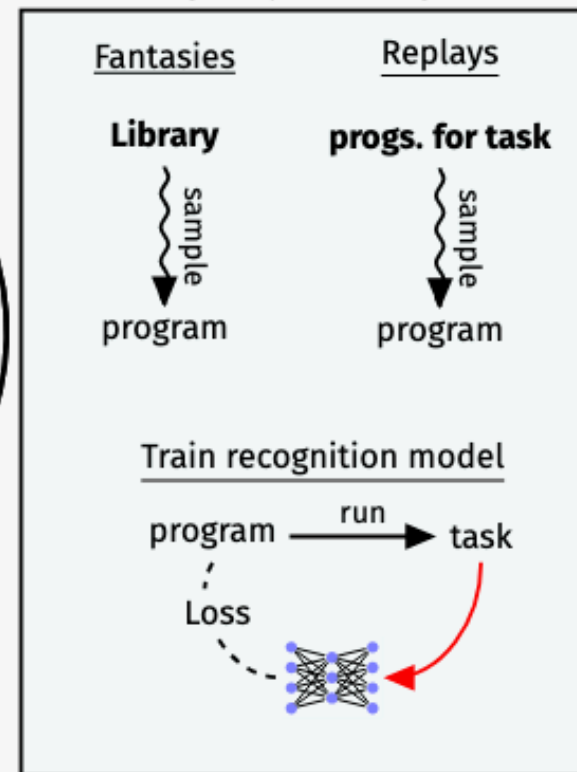
WAKE



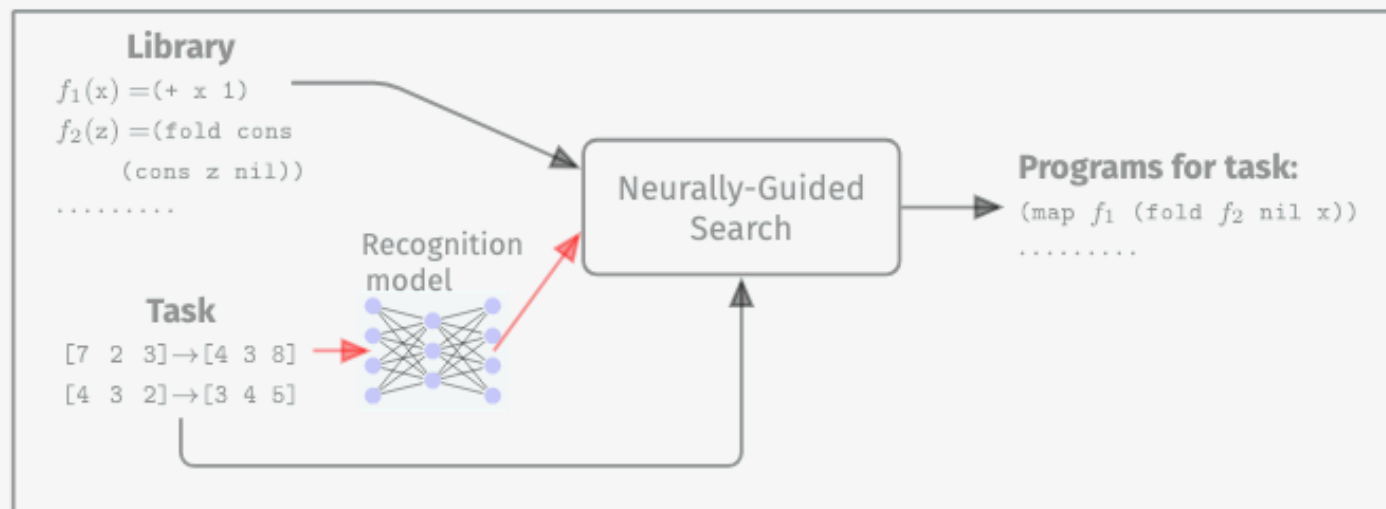
SLEEP: ABSTRACTION



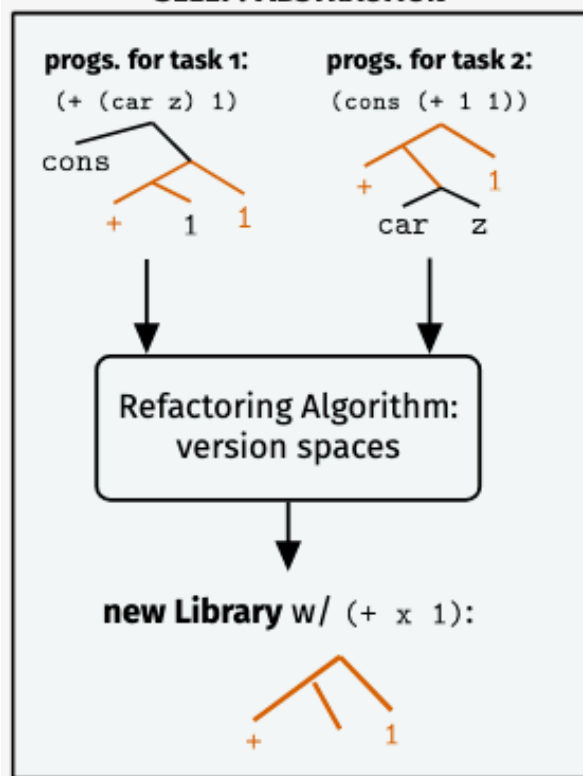
SLEEP: DREAMING



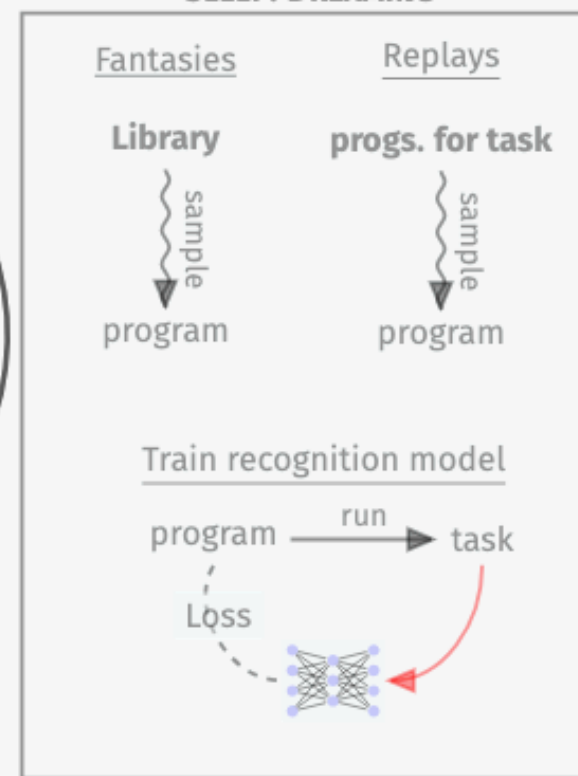
WAKE



SLEEP: ABSTRACTION



SLEEP: DREAMING



Program Induction and learning to learn
learning a DSL
learning to synthesize
synergy between DSL+learned synthesizer

Abstraction Sleep: Growing the library via refactoring

Double

Task: [1 2 3] → [2 4 6]
[4 3 4] → [8 6 8]

Wake: program search

```
(Y (λ (r l) (if (nil? l) nil  
  (cons (+ (car l) (car l))  
    (r (cdr l))))))
```

Decrement

Task: [1 2 3] → [0 1 2]
[4 3 4] → [3 2 3]

Wake: program search

```
(Y (λ (r l) (if (nil? l) nil  
  (cons (- (car l) 1)  
    (r (cdr l))))))
```

Abstraction Sleep: Growing the library via refactoring

Task: [1 2 3] → [2 4 6]
[4 3 4] → [8 6 8]

Wake: program search

```
(Y (λ (r l) (if (nil? l) nil  
              (cons (+ (car l) (car l))
```

Task: [1 2 3] → [0 1 2]
[4 3 4] → [3 2 3]

Wake: program search

```
(Y (λ (r l) (if (nil? l) nil  
              (cons (- (car l) 1)
```

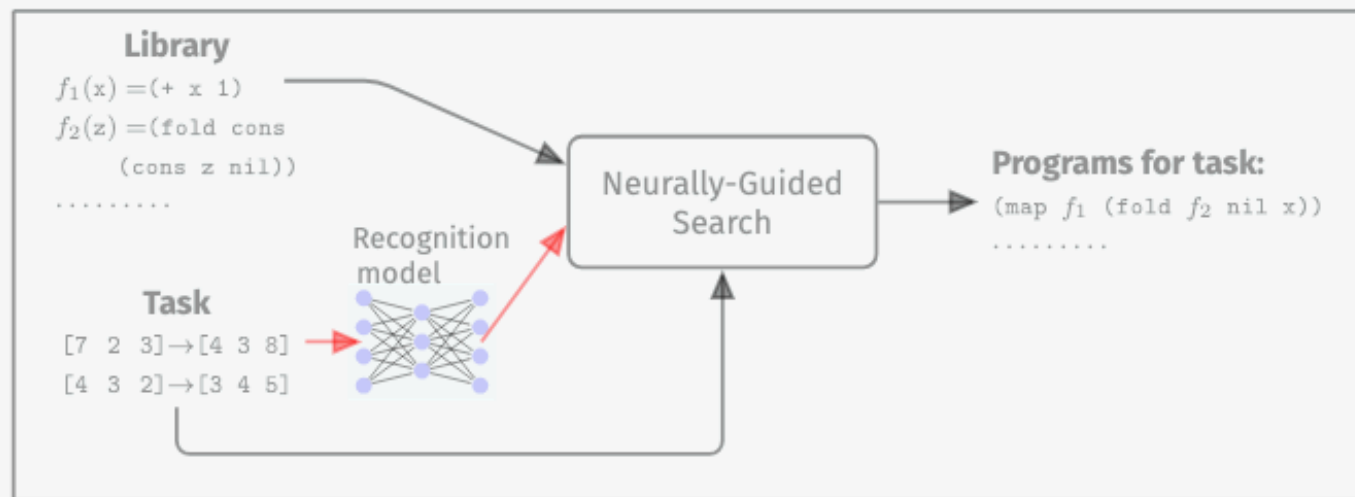
these 10^{14} refactorings represented in exponentially more efficient refactoring data structure: equivalence graphs+version spaces using 10^6 nodes, calculated in under 5min
c.f. [Tate et al 2009], [Gulwani 2012]

(MAP (λ (z) (+ z z))) (MAP (λ (z) (- z 1)))

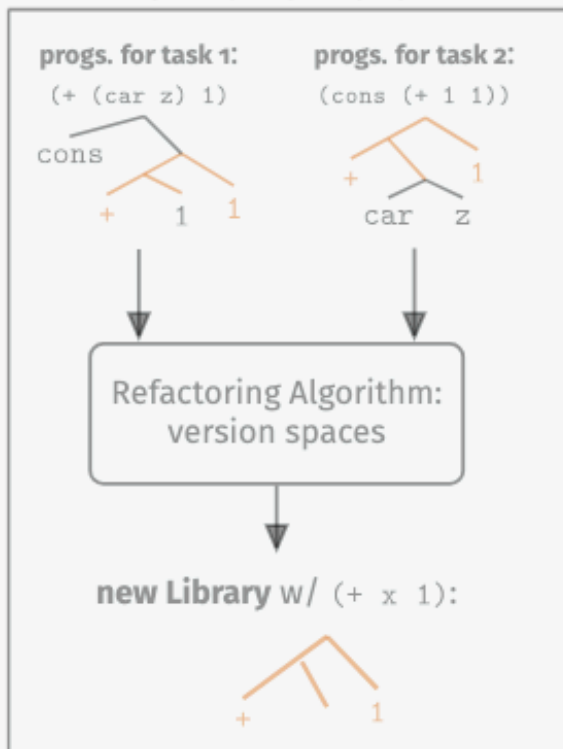
```
MAP = (λ (f) (Y (λ (r l) (if (nil? l) nil  
                          (cons (f (car l))  
                                (r (cdr l)))))))
```

Program Induction and learning to learn
learning a DSL
learning to synthesize
synergy between DSL+learned synthesizer

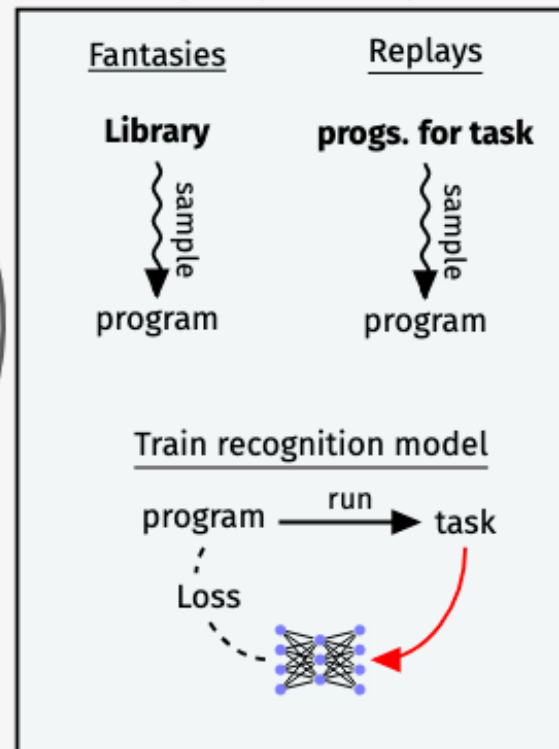
WAKE



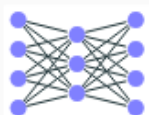
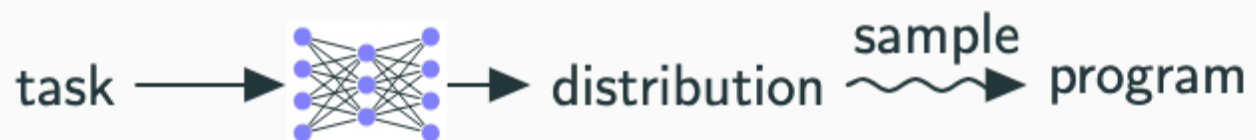
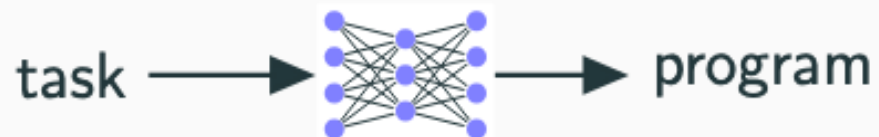
SLEEP: ABSTRACTION



SLEEP: DREAMING



Neural recognition model guides search

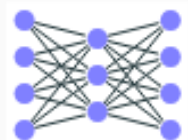
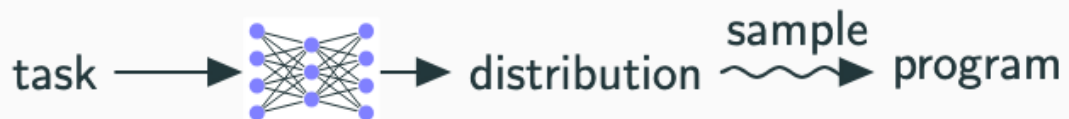
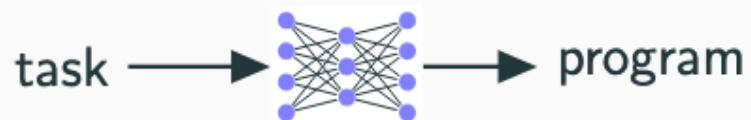


is a...

recurrent network (Devlin et al 2017)

unigram model (Menon et al 2013; Balog et al 2016)

Neural recognition model guides search



is a **“bigram”** model over syntax trees



Program Induction and learning to learn
learning a DSL
learning to synthesize
synergy between DSL+learned synthesizer

DreamCoder Domains

List Processing

Sum List

[1 2 3] → 6

[4 6 8 1] → 17

Double

[1 2 3] → [2 4 6]

[4 5 1] → [8 10 2]

Text Editing

Abbreviate

Allen Newell → A.N.

Herb Simon → H.S.

Drop Last Three

shrdlu → shr

shakey → sha

Regexes

Phone numbers

(555) 867-5309

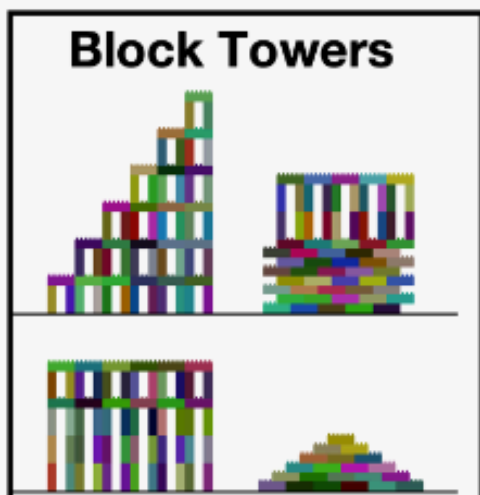
(650) 555-2368

Currency

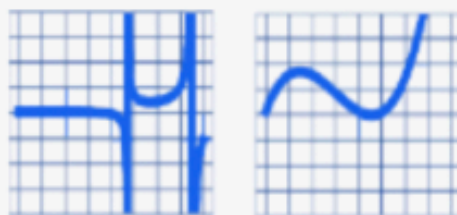
\$100.25

\$4.50

LOGO Graphics



Symbolic Regression



$$y = f(x)$$

Recursive Programming

Filter Red

[red red blue blue] → [blue blue]

[red black green red black green] → [black green black green]

[red black black red] → [black black]

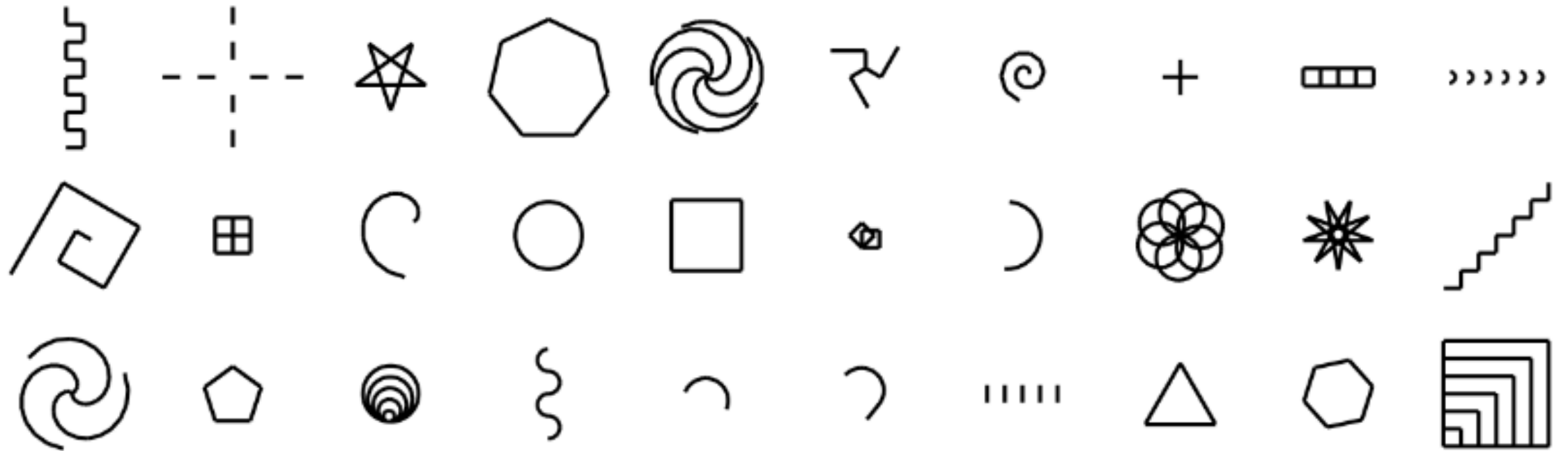
Physical Laws

$$\vec{a} = \frac{1}{m} \sum_i \vec{F}_i$$

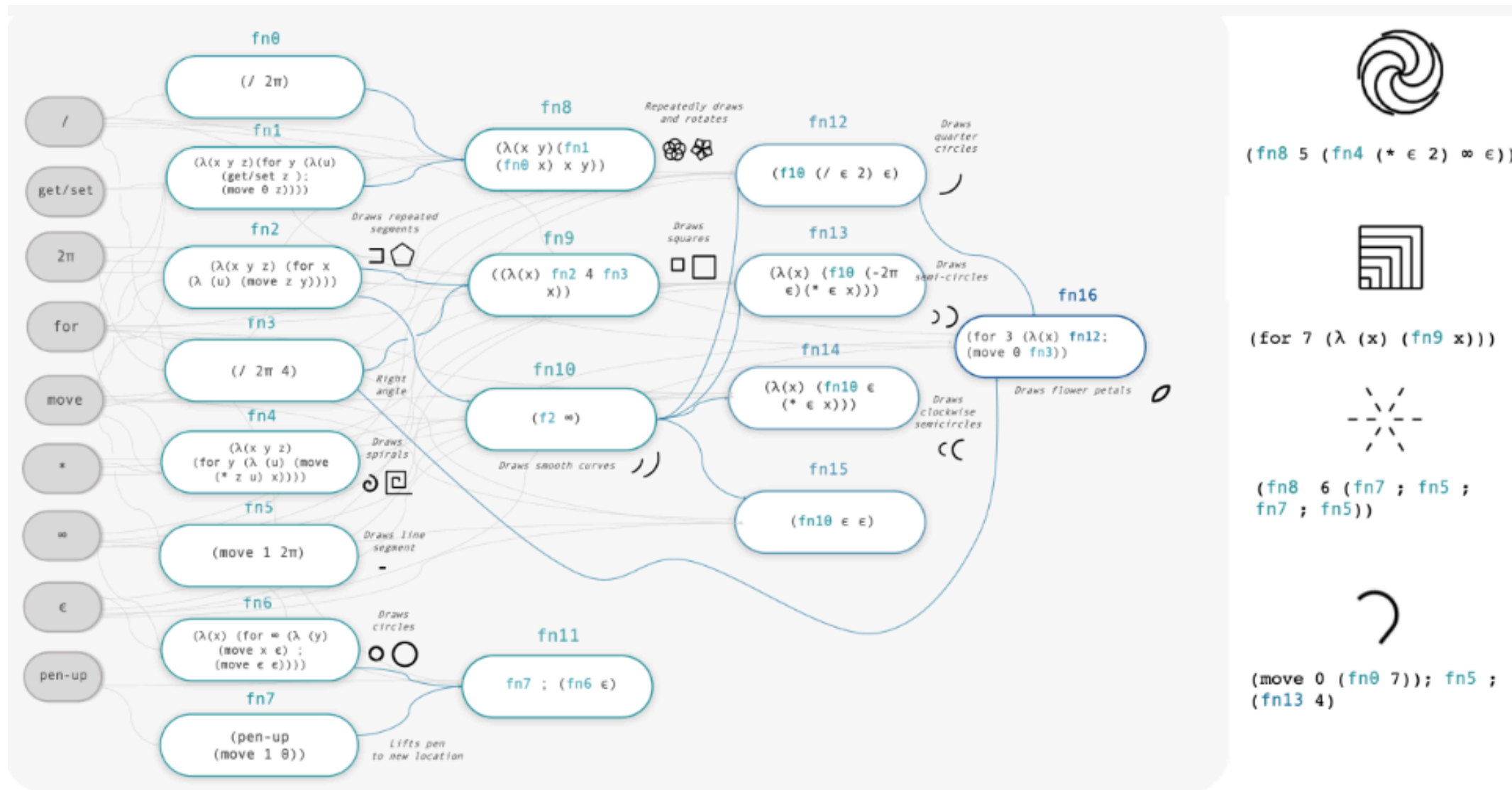
$$\vec{F} \propto \frac{q_1 q_2}{|\vec{r}|^2} \hat{r}$$

LOGO Turtle Graphics

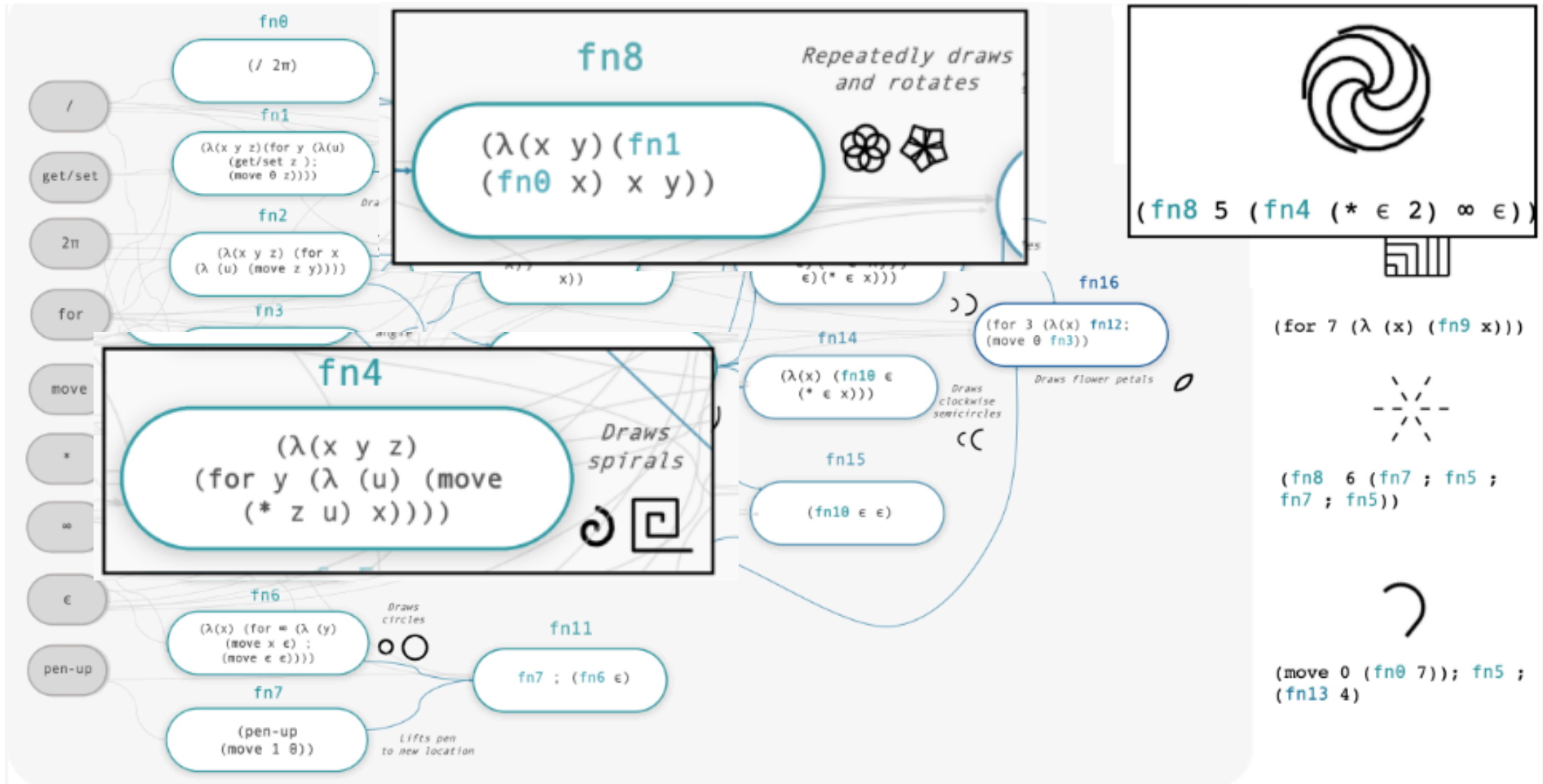
30 out of 160 tasks



LOGO Turtle Graphics – learning an interpretable library



LOGO Turtle Graphics – learning an interpretable library



LOGO Turtle Graphics – learning an interpretable library

fn0
(/ 2π)

fn1
(λ(x y z)(for y (λ(u) (get/set z) : (move θ z))))

fn2
(λ(x y z) (for x (λ (u) (move z y))))

fn3
(/ 2π 4)

fn4
(λ(x y z) (for y (λ (u) (move (* z u) x))))

fn5
(move 1 2π)

fn6
(λ(x) (for = (λ (y) (move x ε) : (move ε ε))))

fn7
(pen-up (move 1 θ))

radial symmetry(n, body)

fn16
(λ(x) fn12: (3))
flower petals

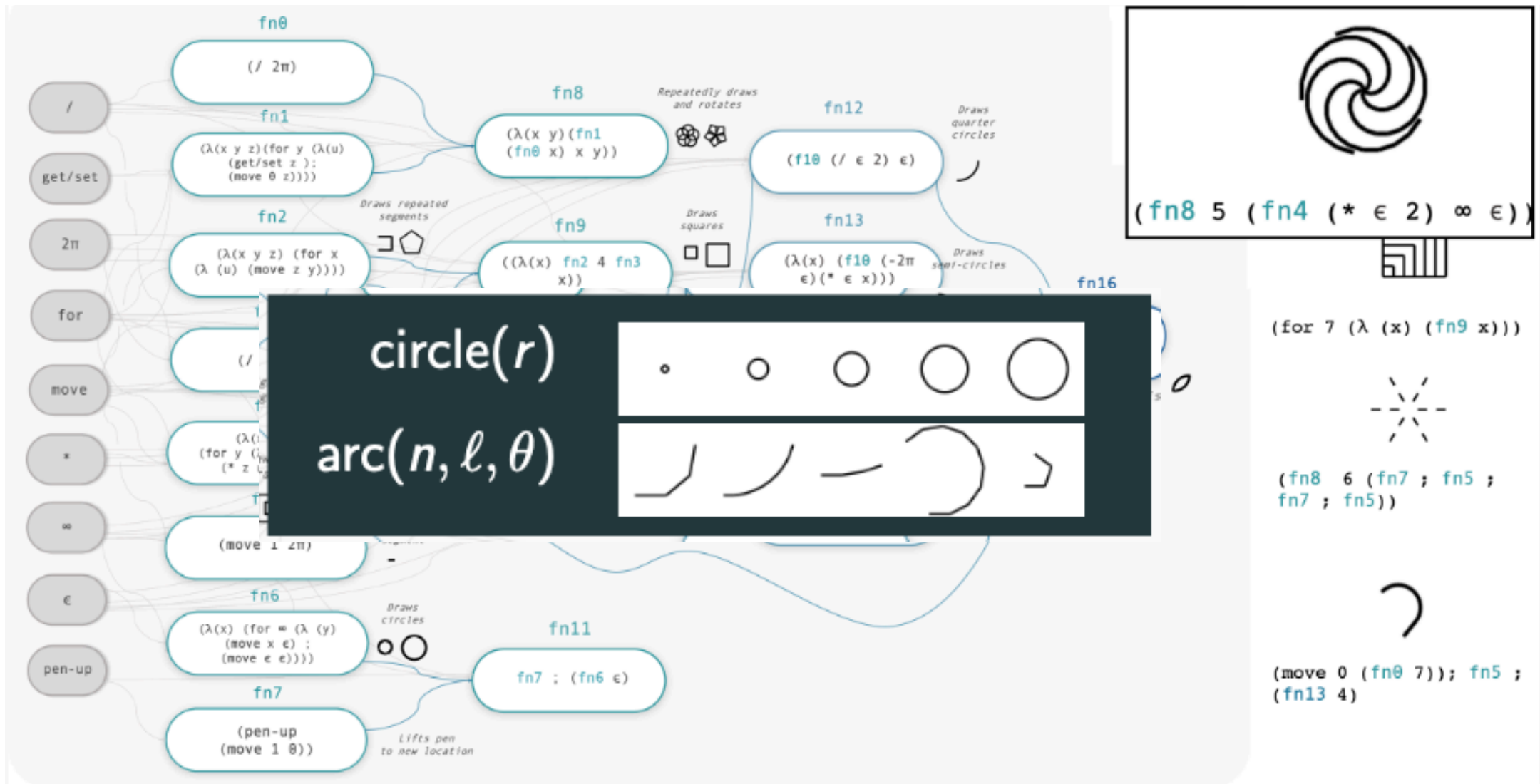
```
(fn8 5 (fn4 (* ε 2) ∞ ε))
```

```
(for 7 (λ (x) (fn9 x)))
```

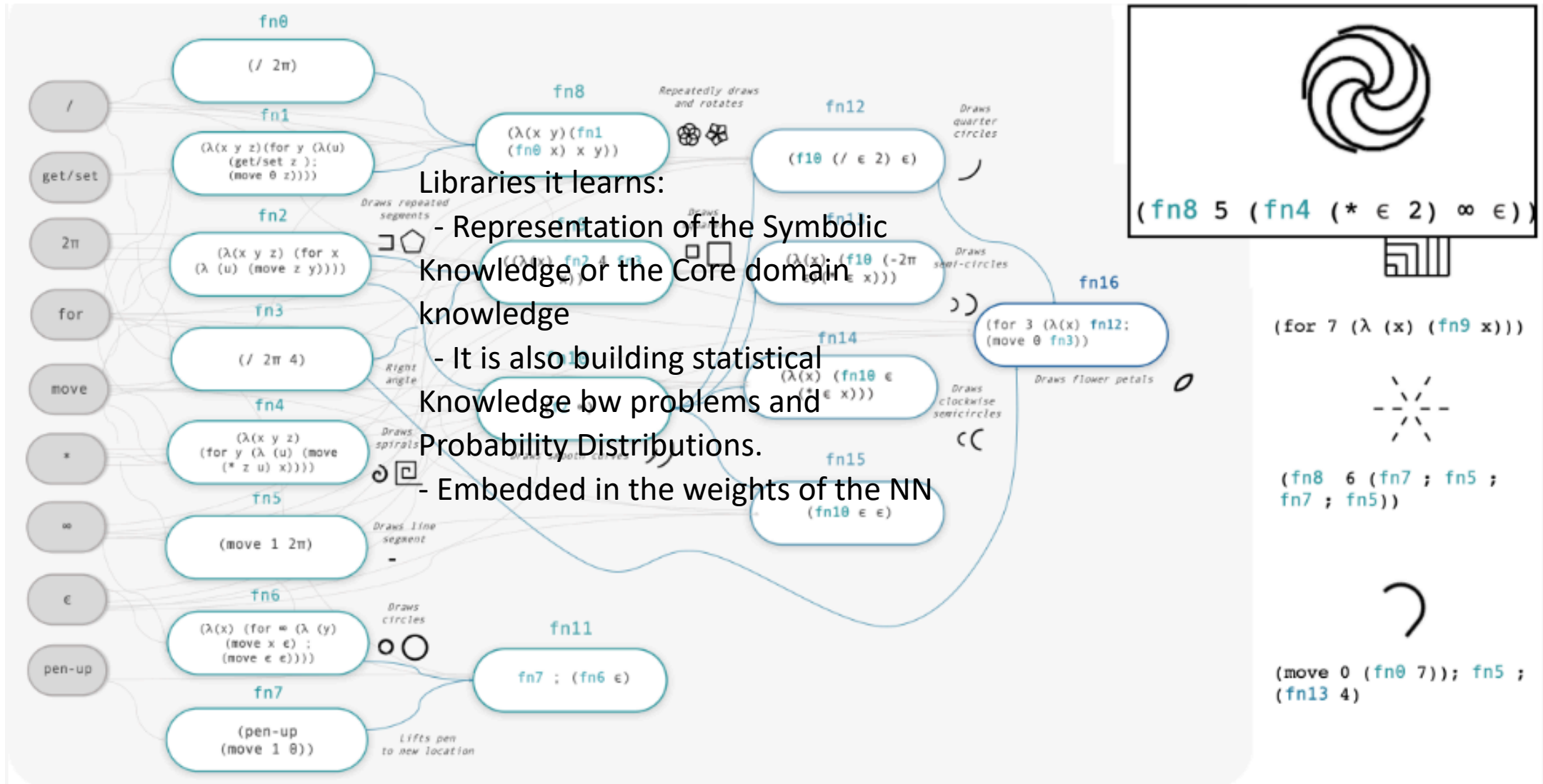
```
(fn8 6 (fn7 ; fn5 ; fn7 ; fn5))
```

```
(move 0 (fn0 7)); fn5 ; (fn13 4)
```

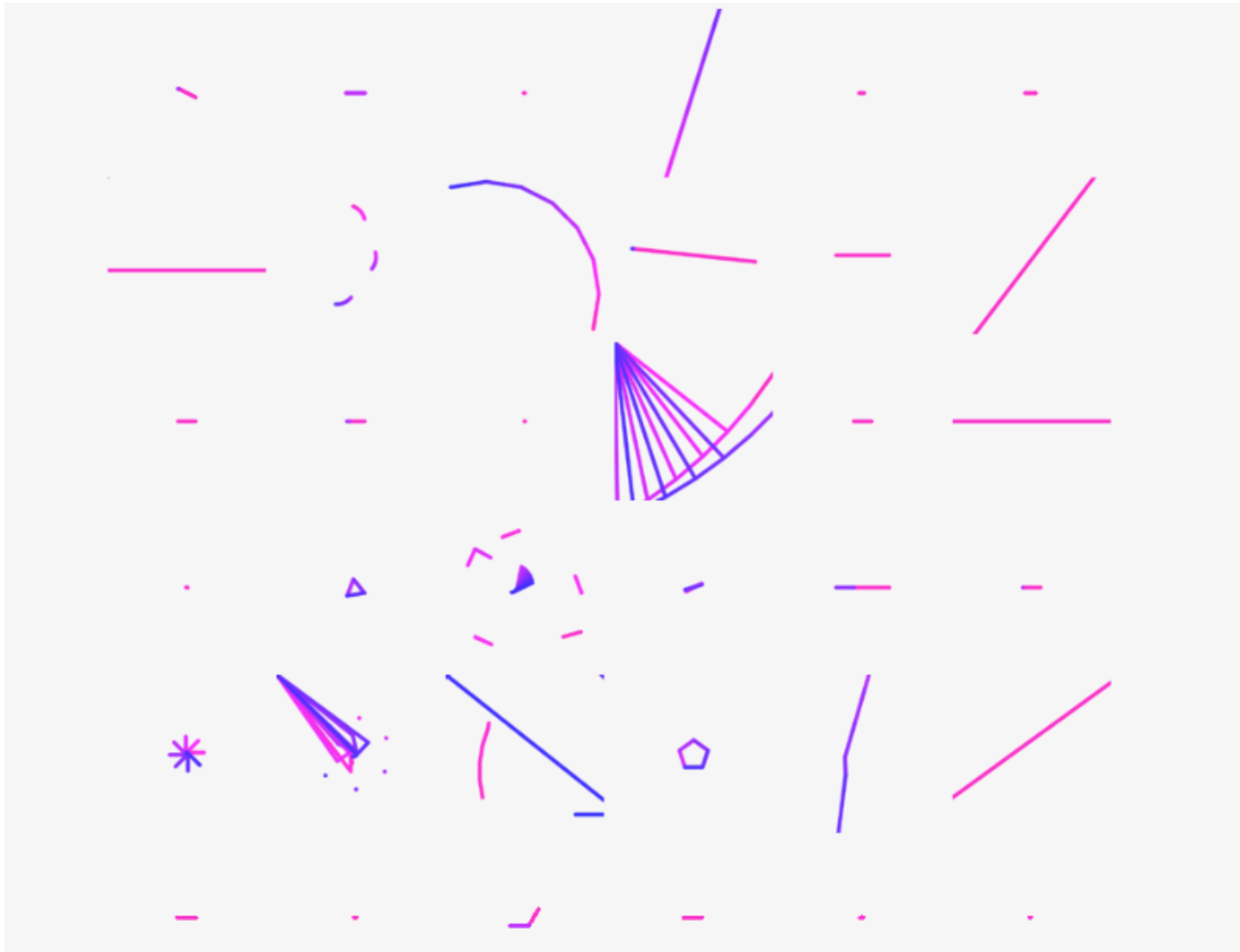
LOGO Turtle Graphics – learning an interpretable library



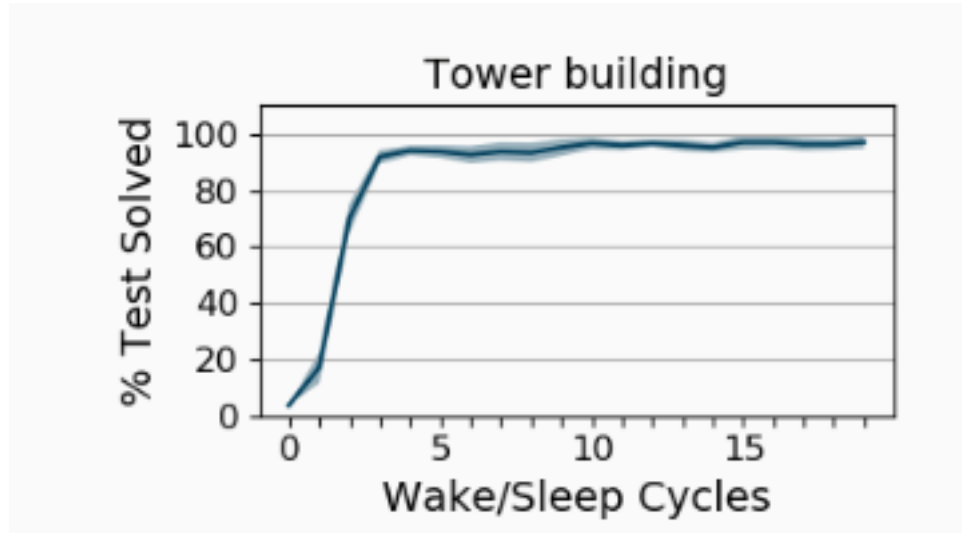
LOGO Turtle Graphics – learning an interpretable library



What does DreamCoder dream of? (before learning)

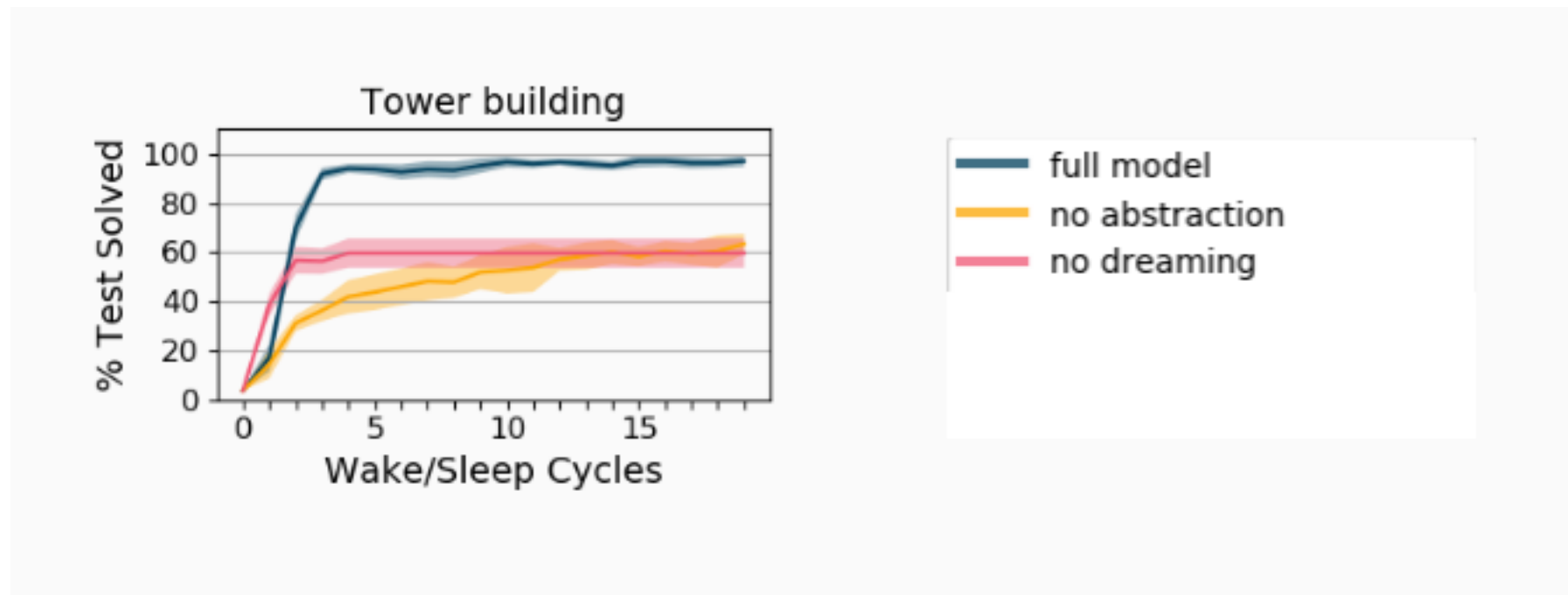


Learning dynamics

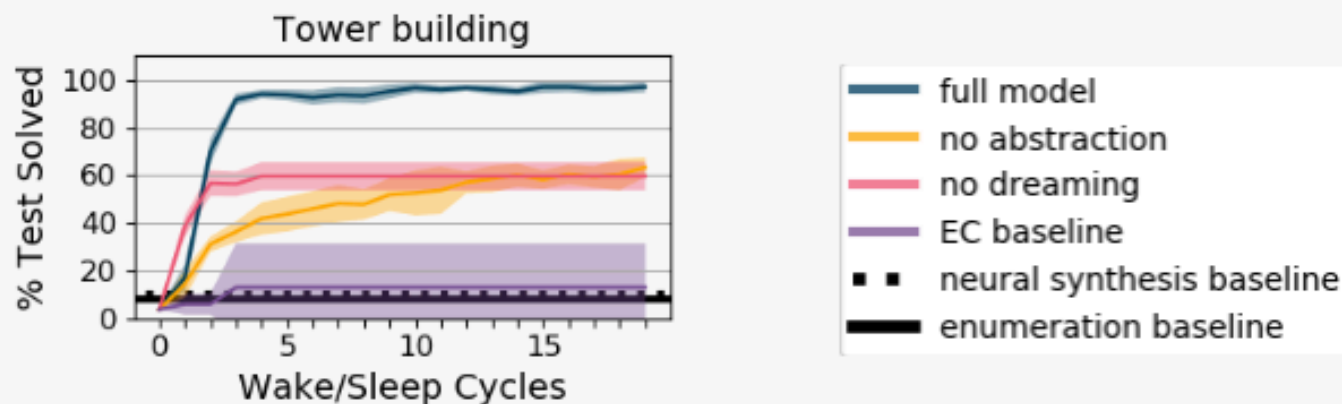


Bootstrapping action

Learning dynamics

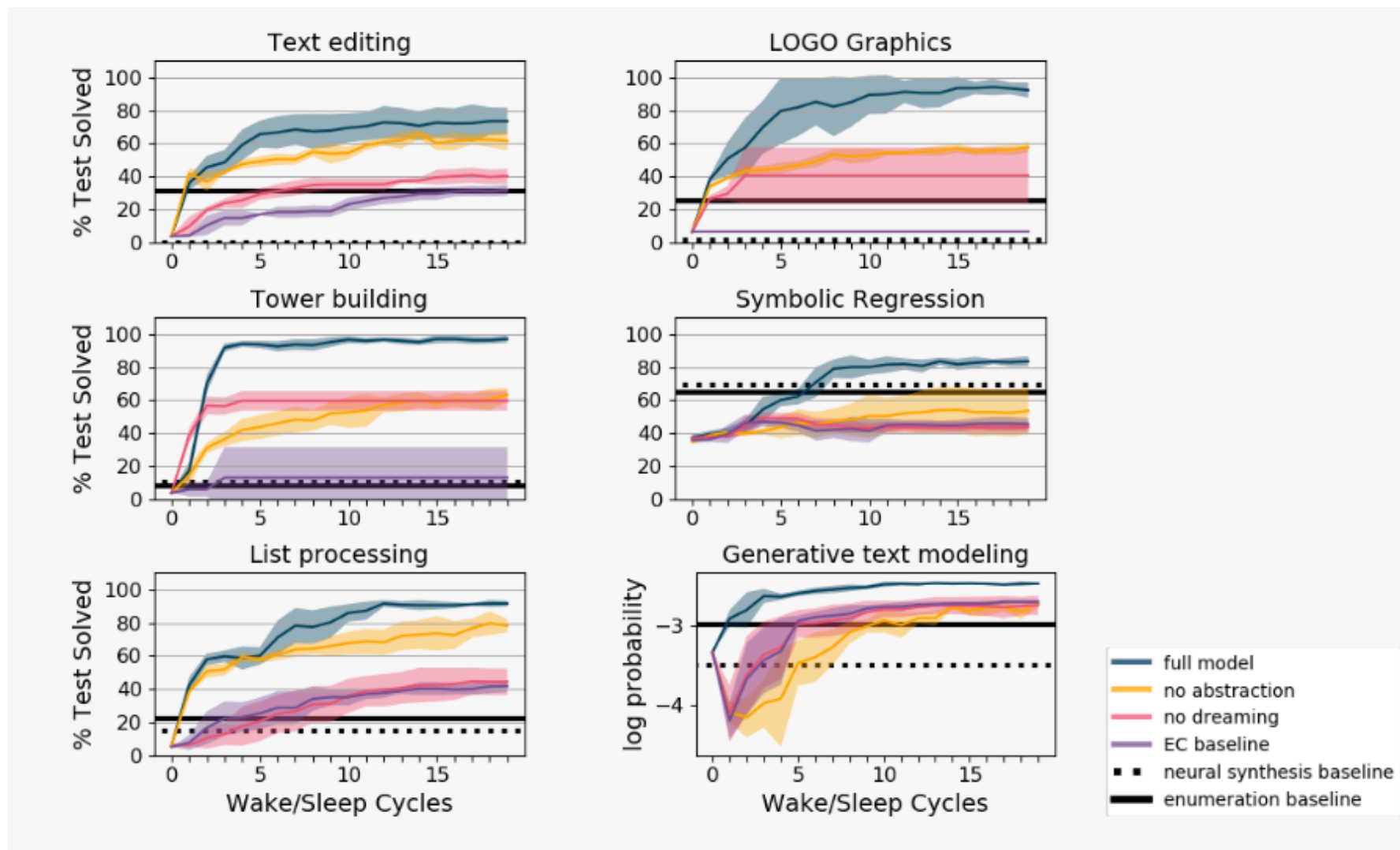


Learning dynamics

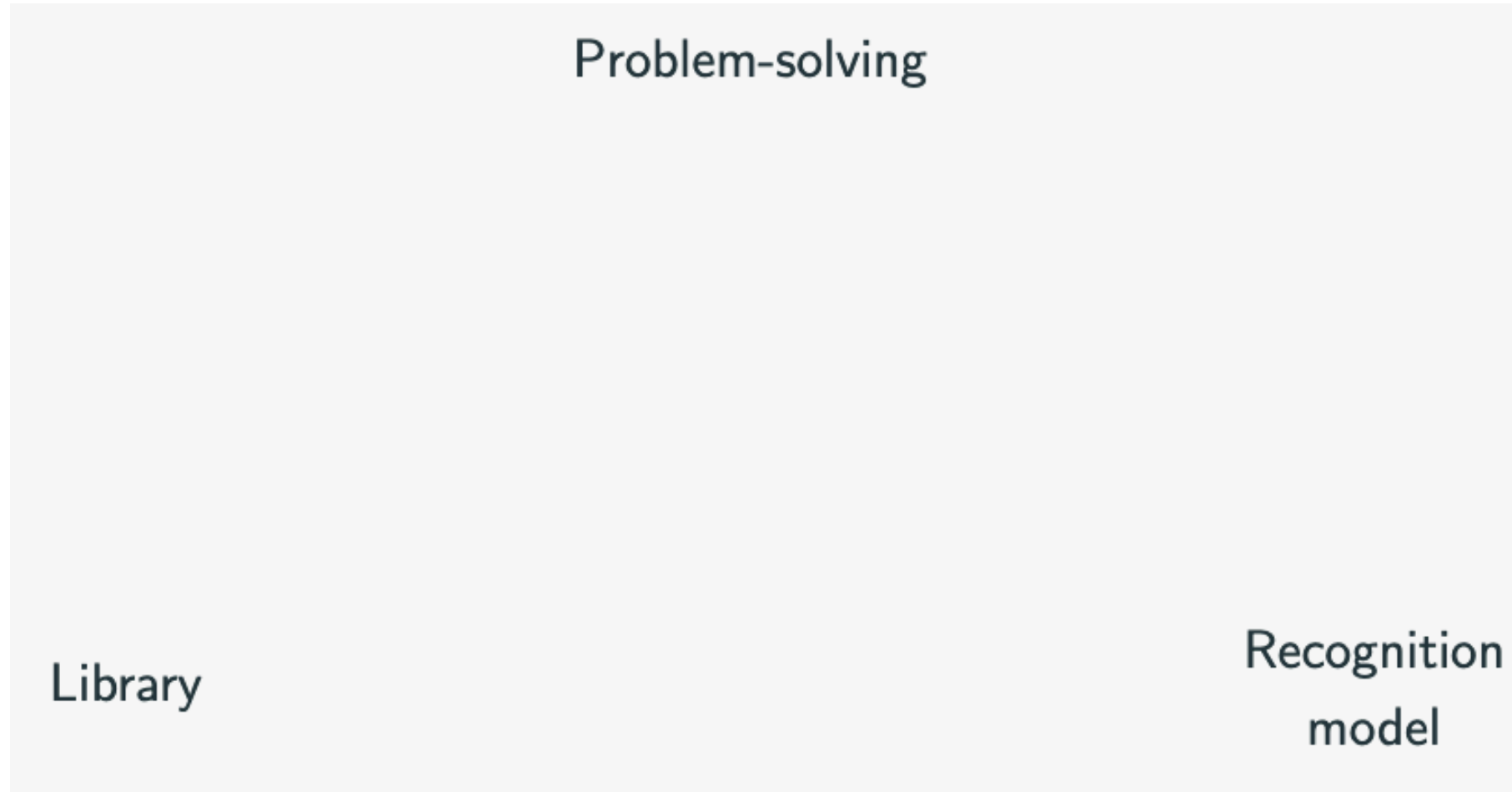


baselines: Exploration-Compression, EC [Dechter et al. 2013]
neural program synthesis, RobustFill [Devlin et al. 2017]
24 hours of brute-force enumeration

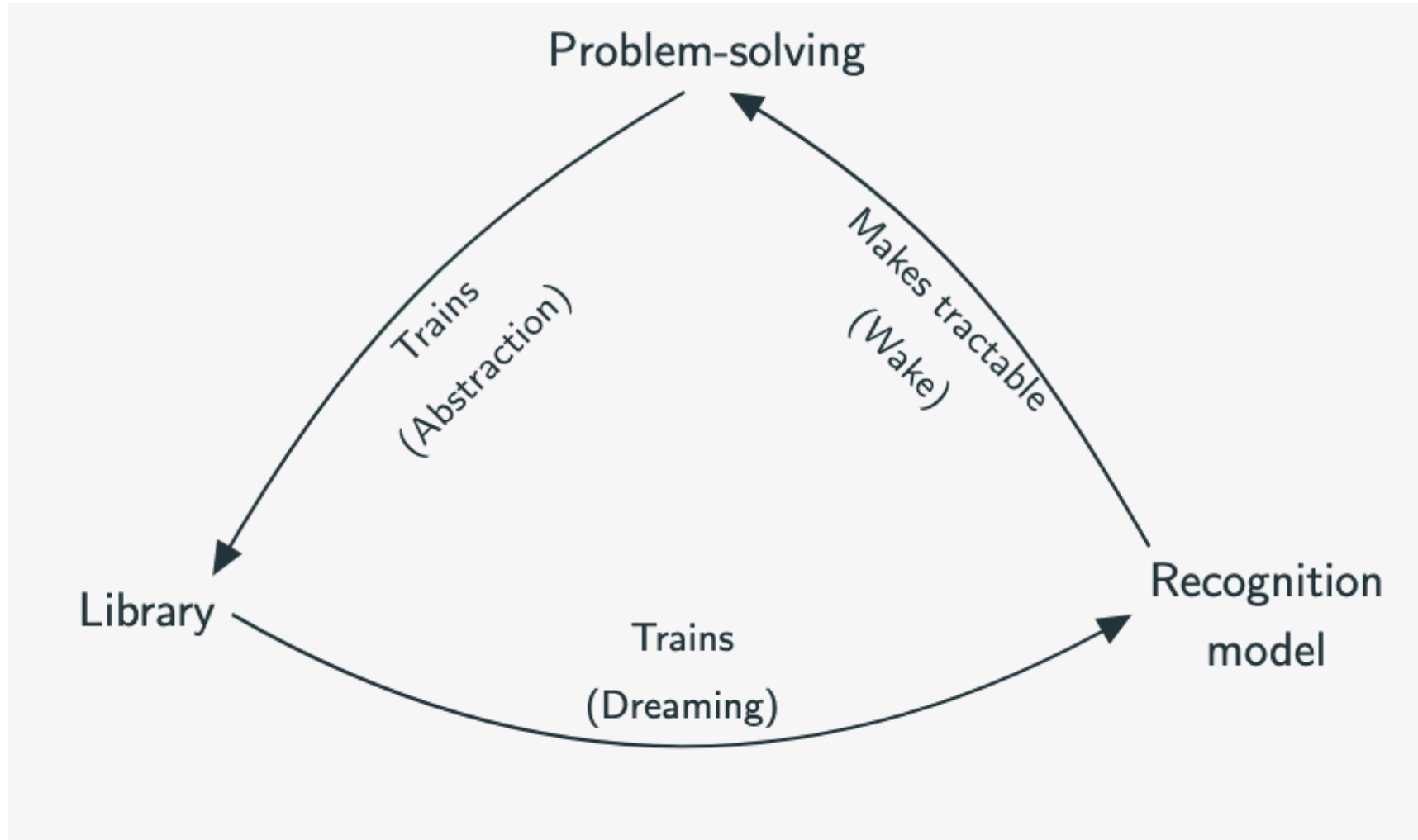
Learning dynamics



Synergy between recognition model and library learning



Synergy between recognition model and library learning



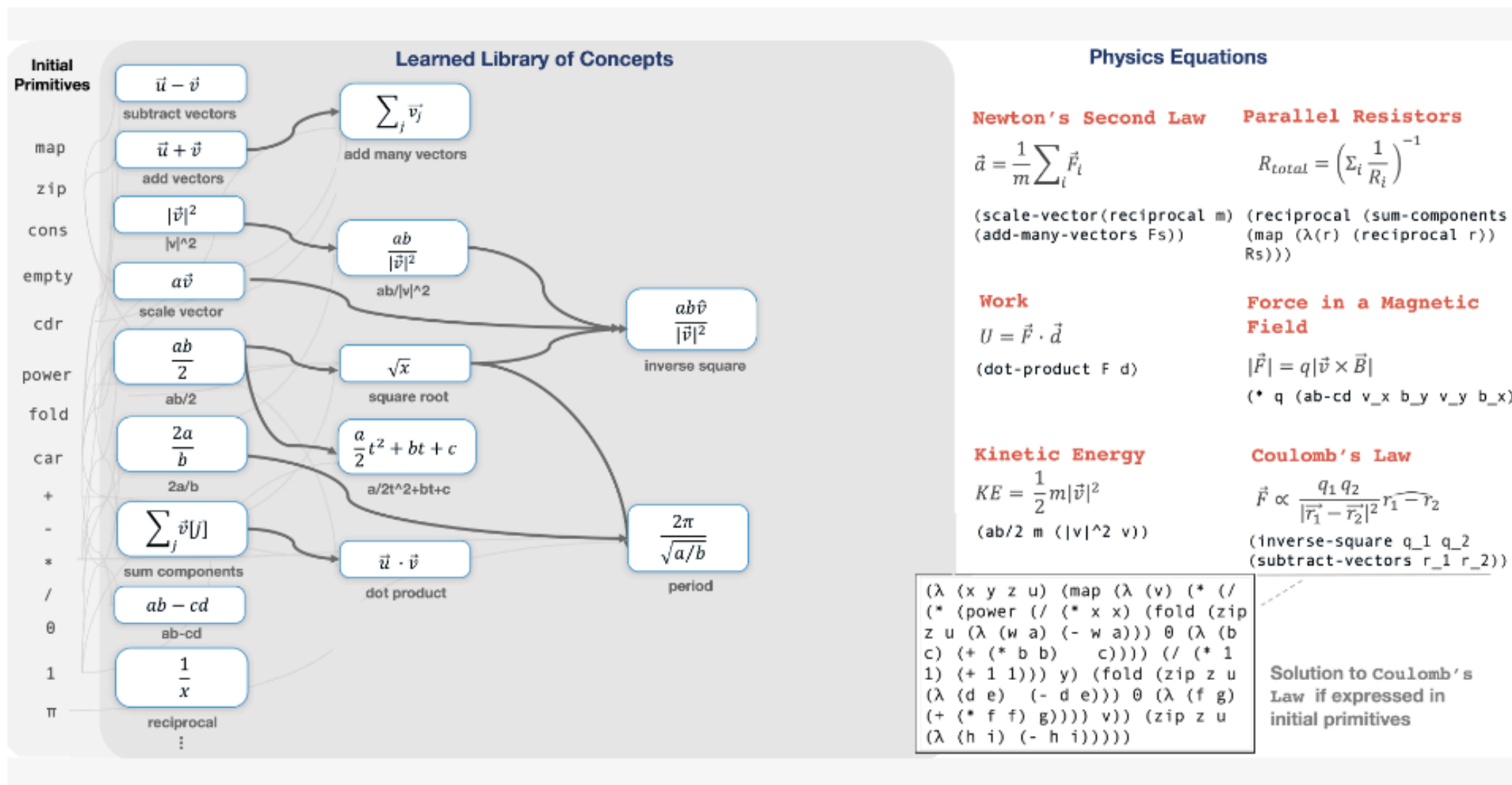
From learning libraries, to learning languages

modern functional programming → physics

From learning libraries,
to learning languages

1950's Lisp → modern functional programming → physics

Growing languages for vector algebra and physics



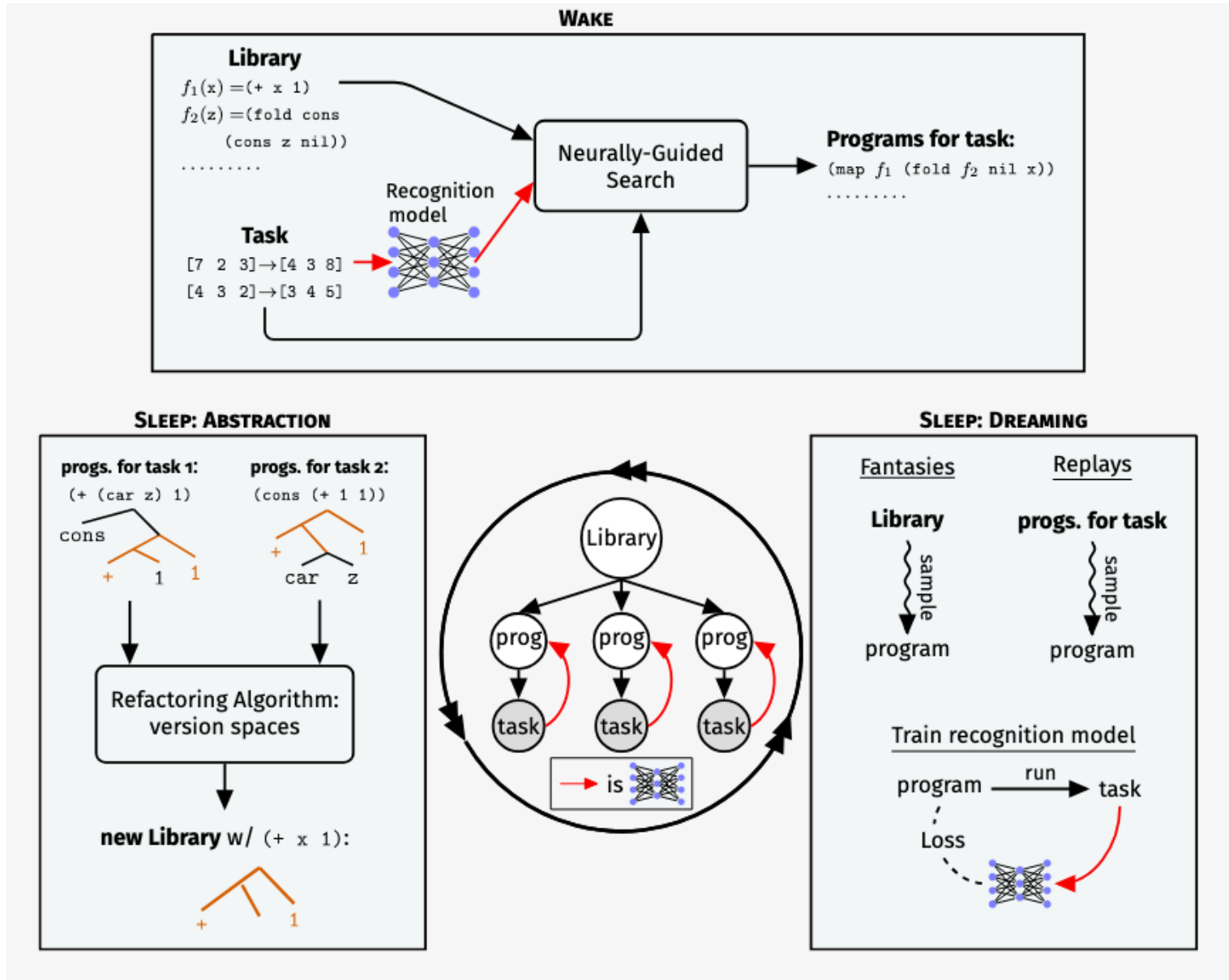
Lessons

Library learning interacts synergistically with neural synthesis:
bootstrapping, more than sum of parts

Symbols aren't necessarily interpretable. Grow the language based on experience to make it more powerful *and* more human understandable

Learning-from-scratch is possible in principle. Don't do it. But program induction makes it convenient to build in what we know how to build in, and then learn on top of that

end.



Logistical

- Next week mid-progress meeting for the projects.
- Re-scheduling Fridays class.
- New-paper