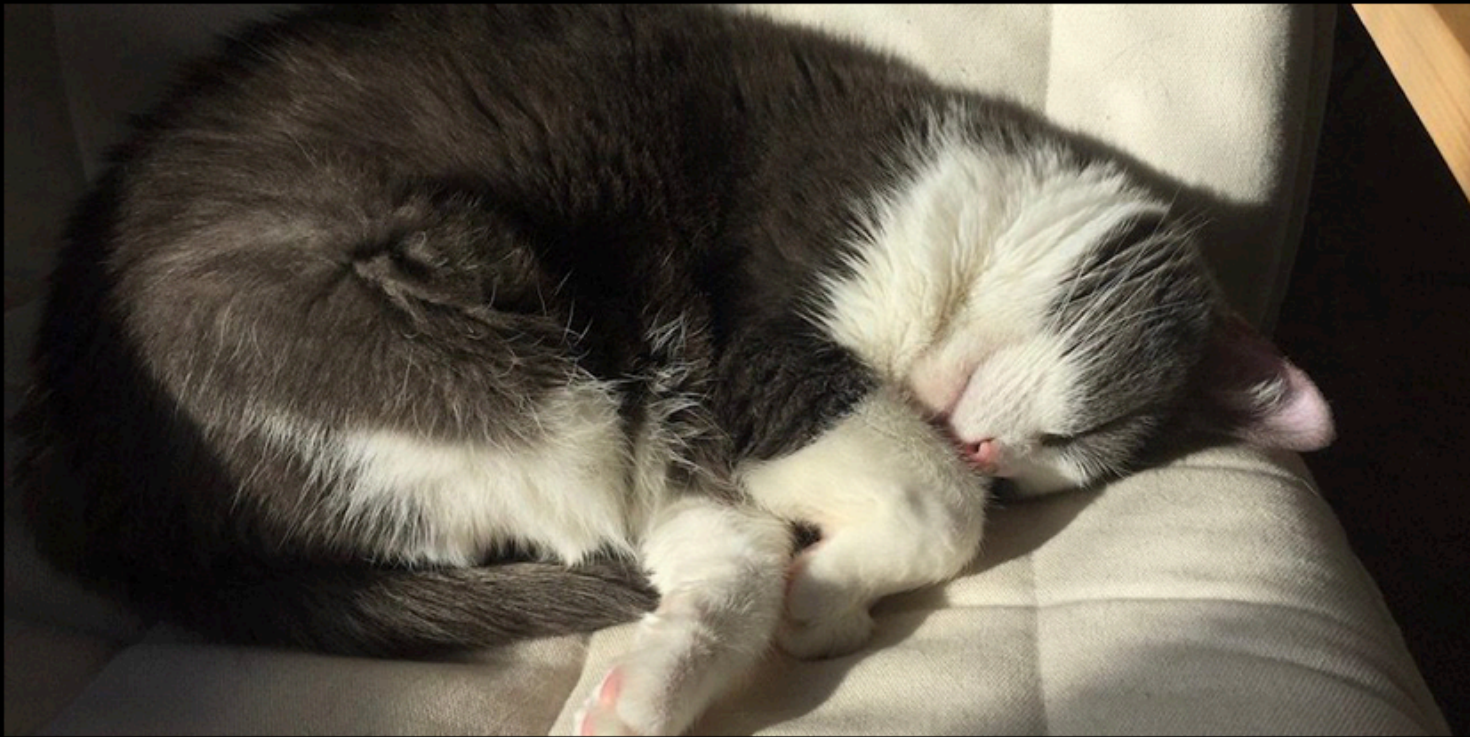


# ROBUSTNESS AS A REFINEMENT TYPE

Wen Kokke, Ekaterina Komendantskaya, and Daniel Kienitz

Lab for AI and Verification, Heriot-Watt University



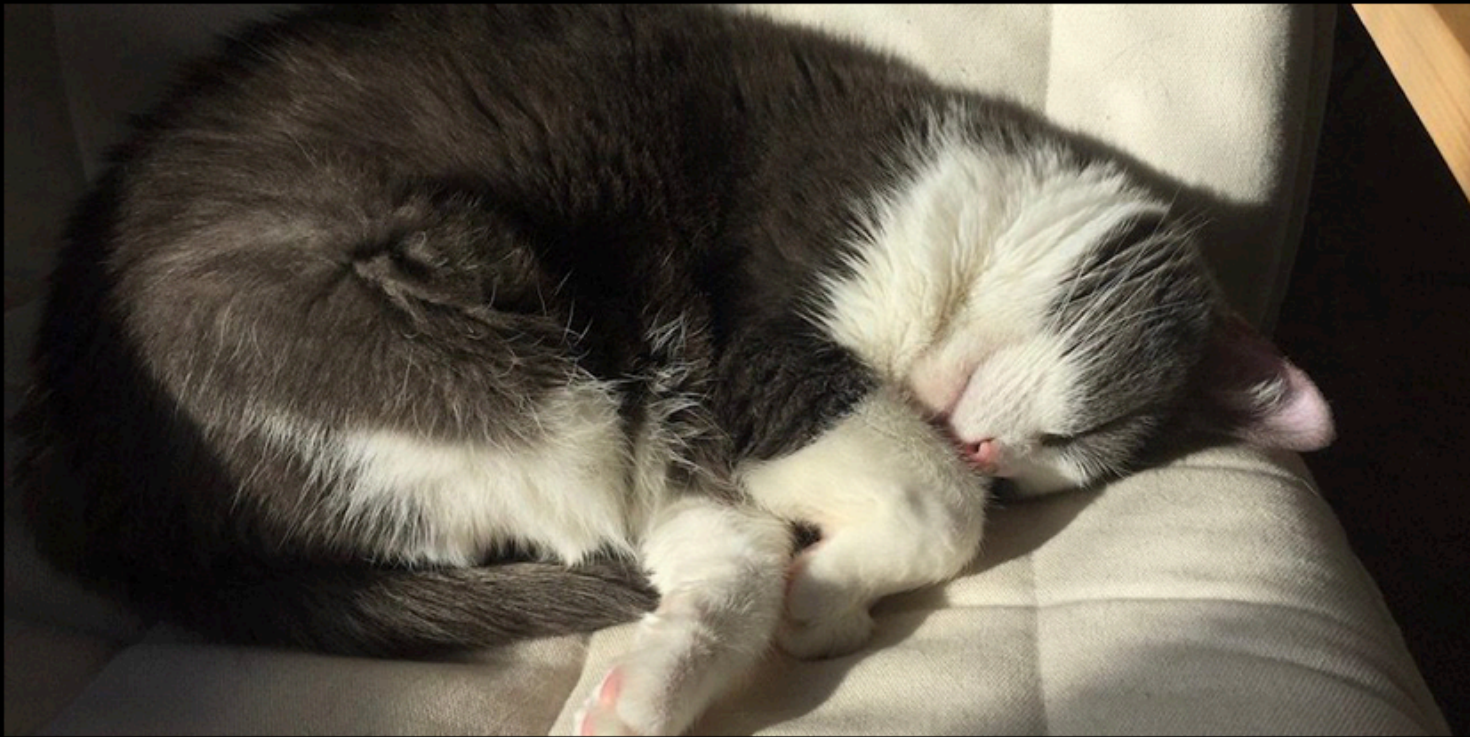




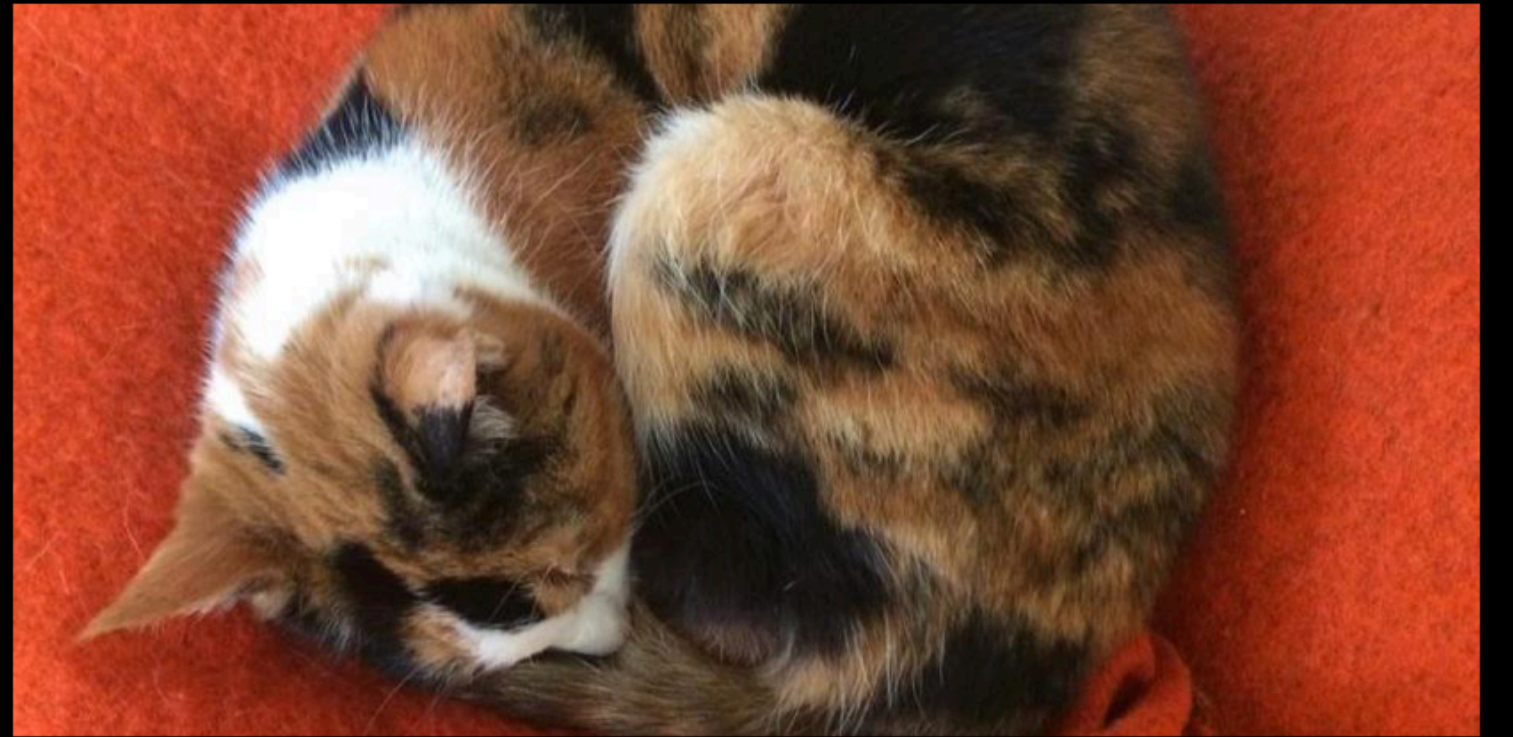
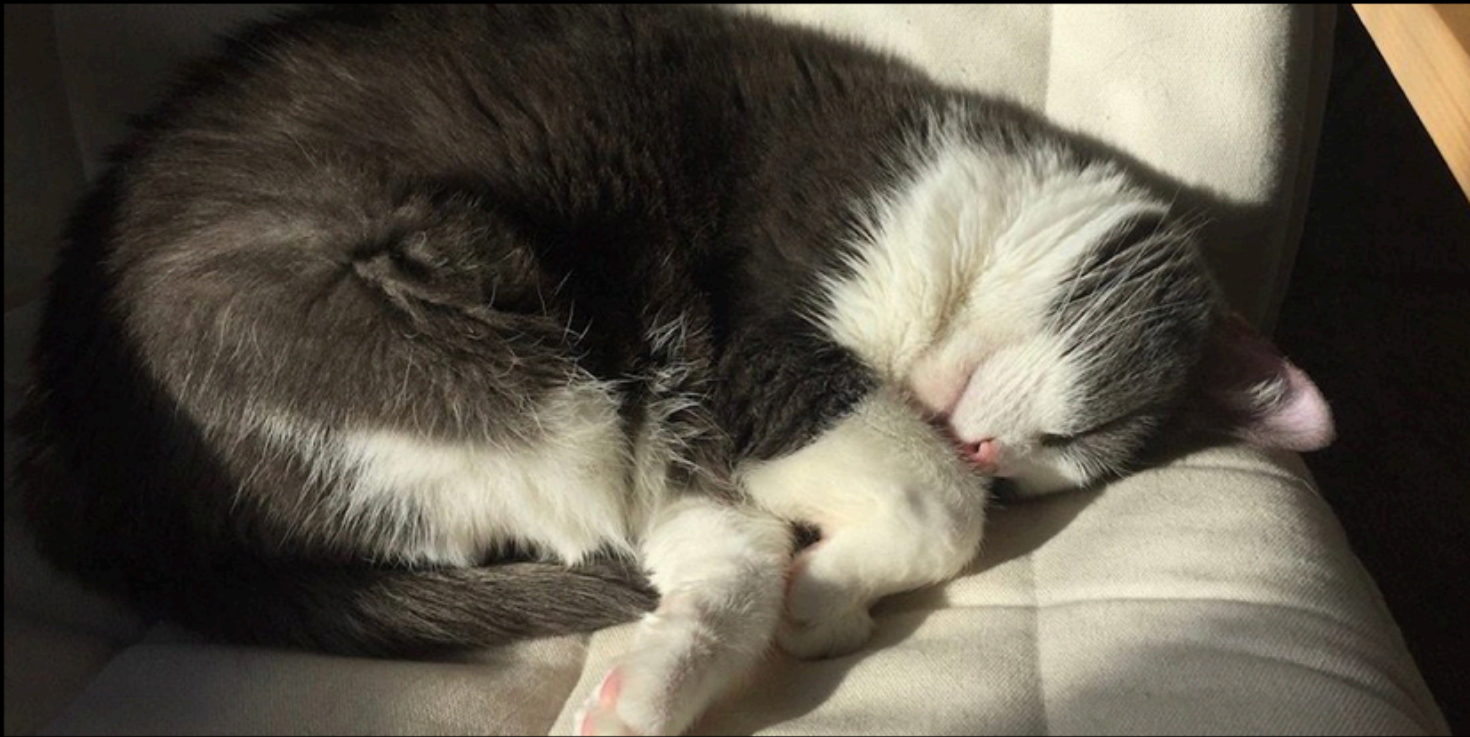




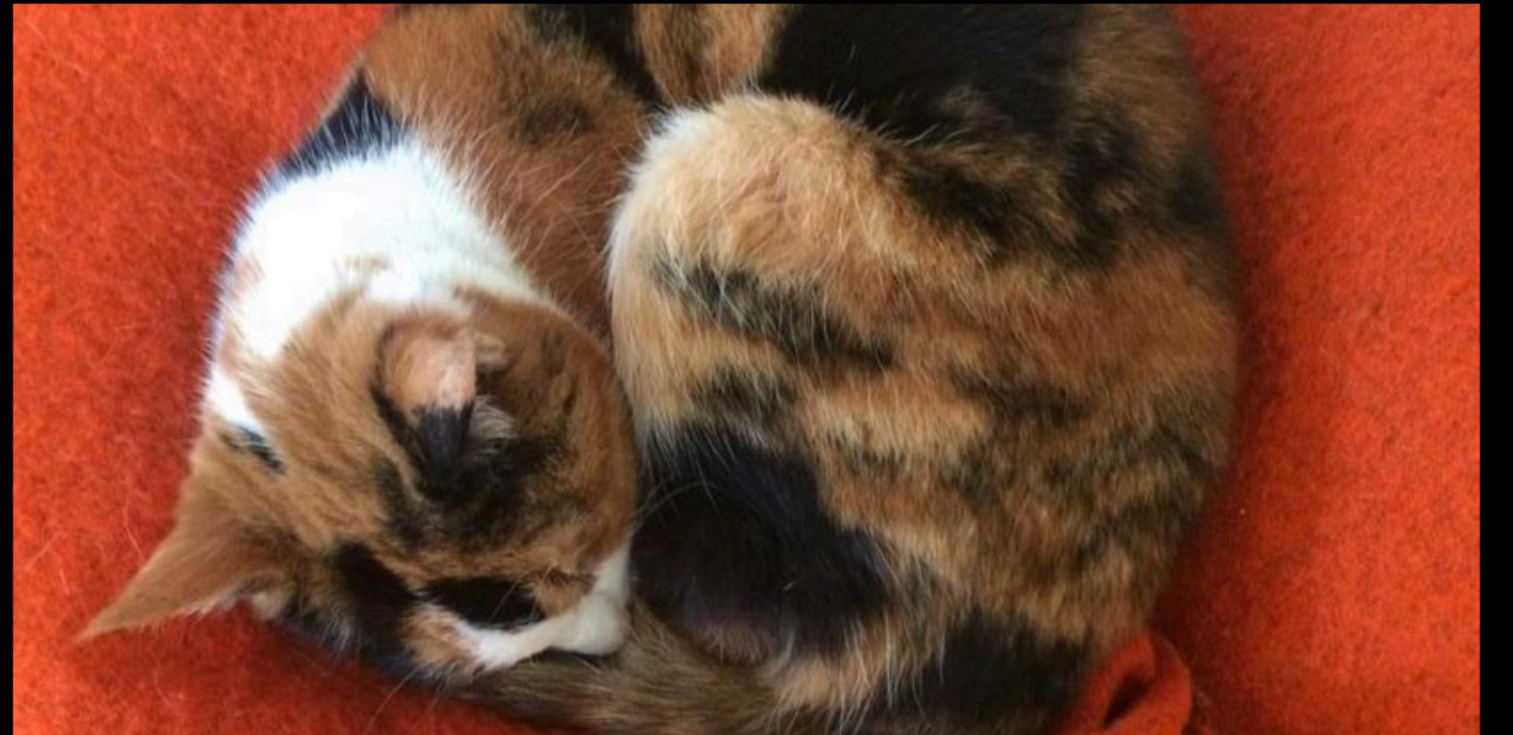
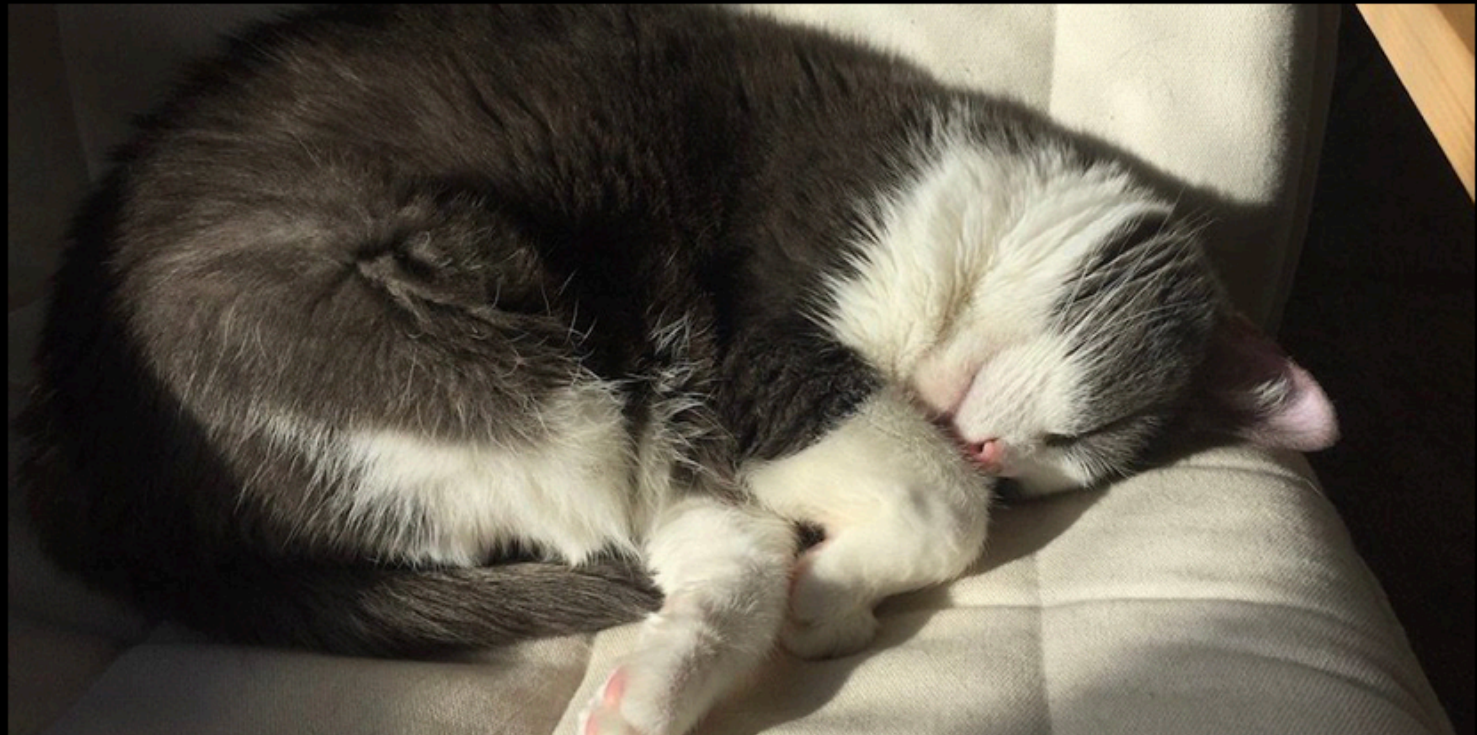








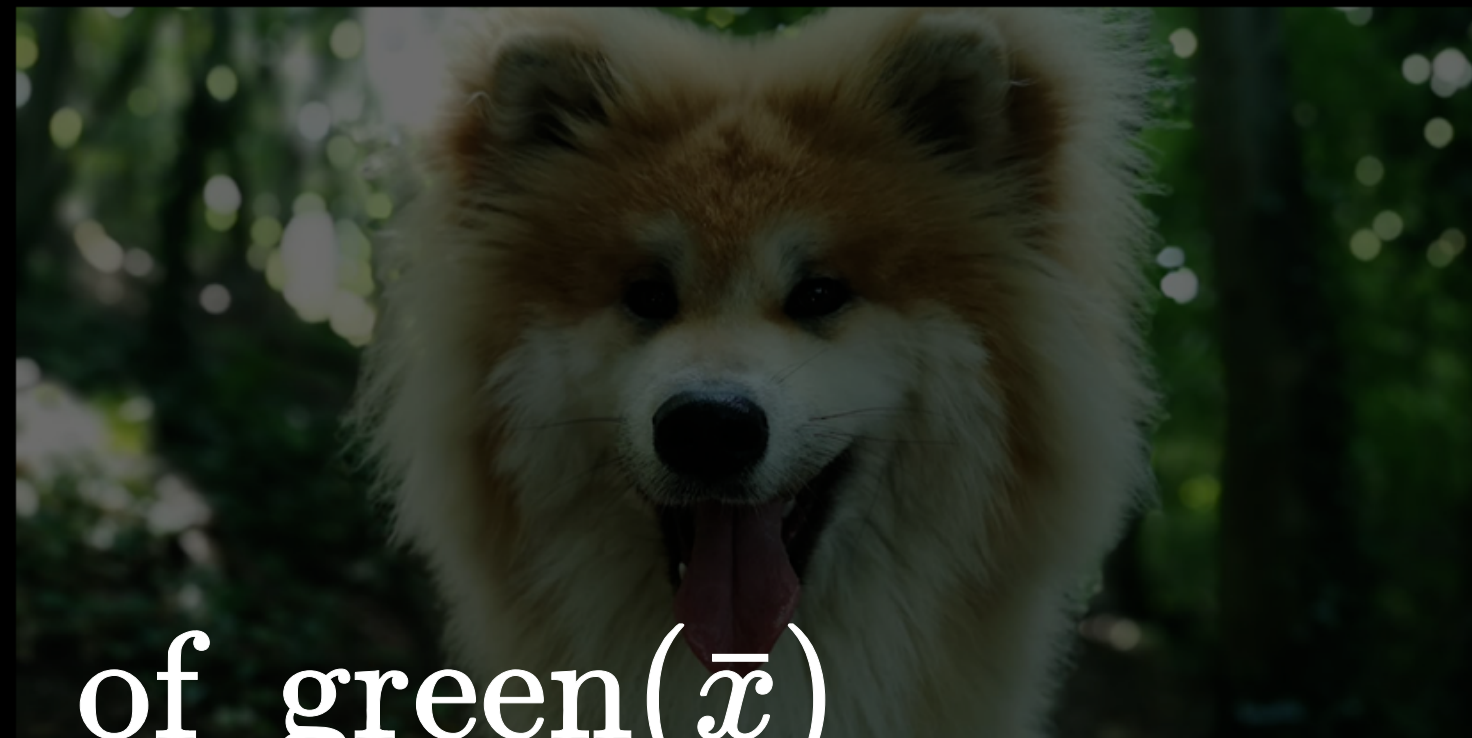
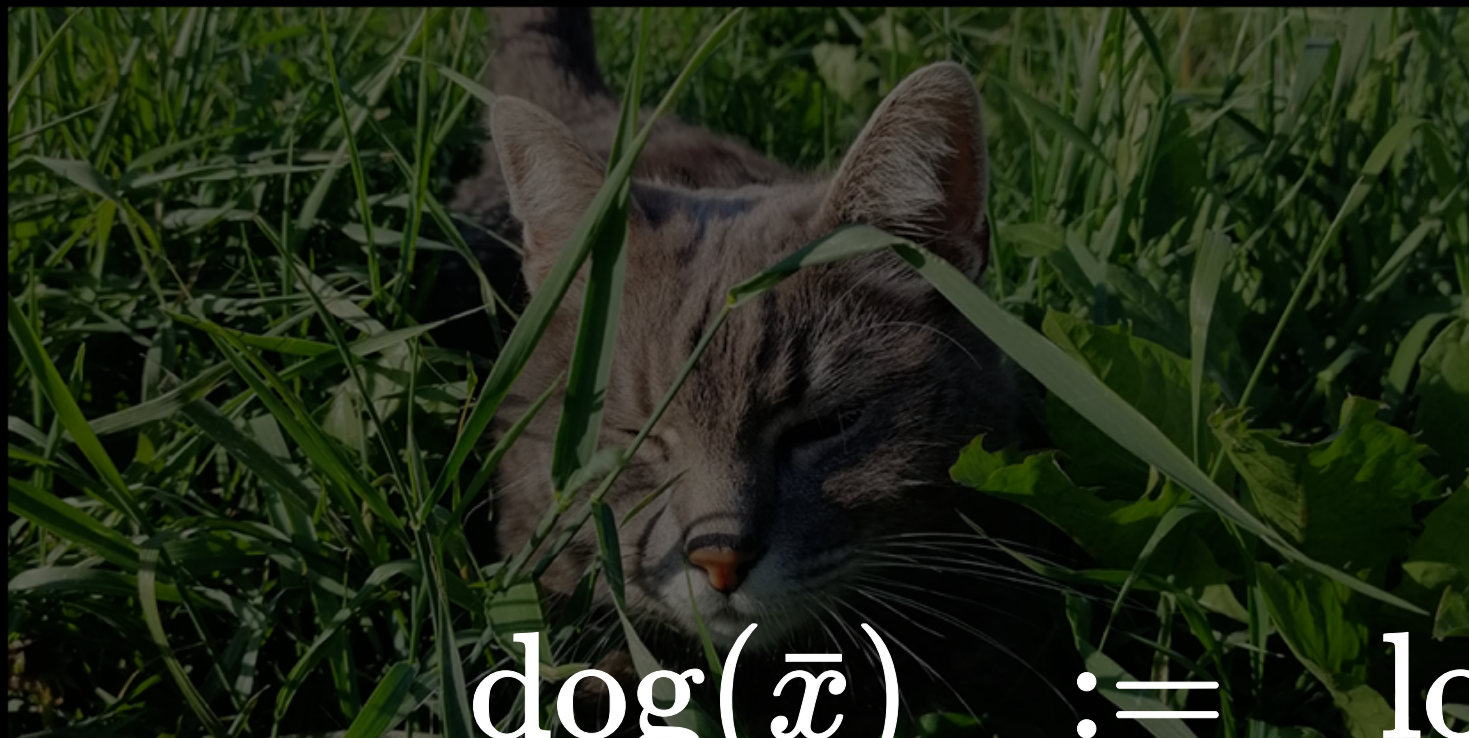










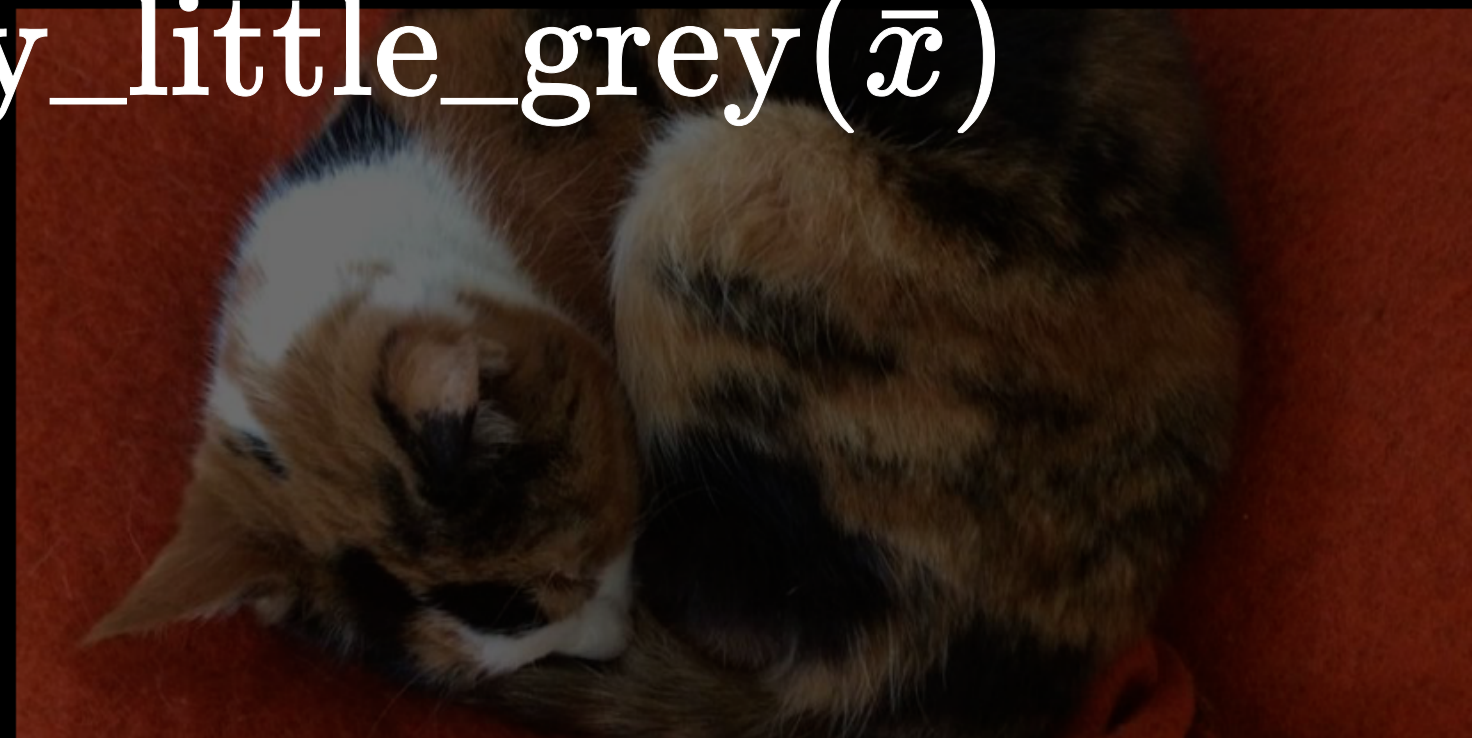


$\text{dog}(\bar{x}) \quad := \quad \text{lots\_of\_green}(\bar{x})$



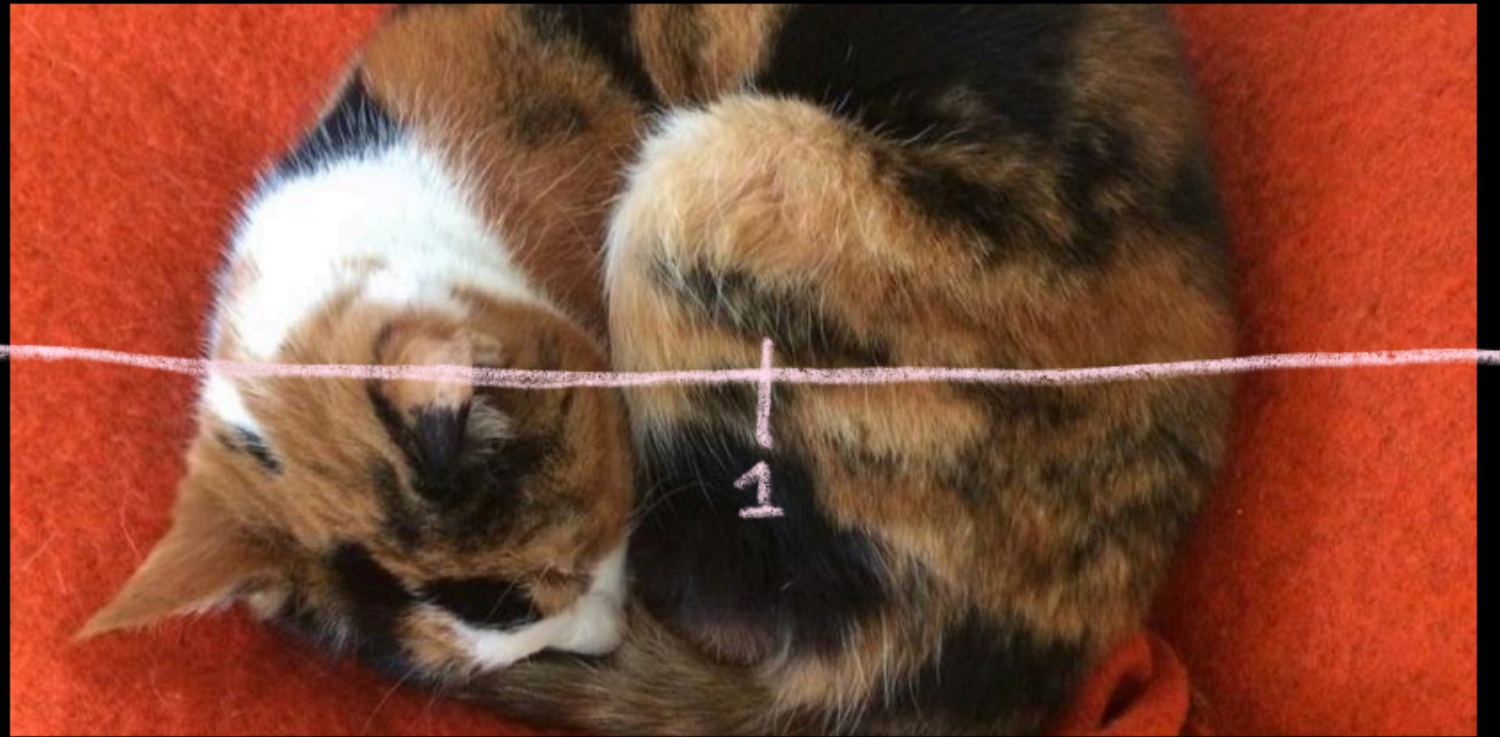
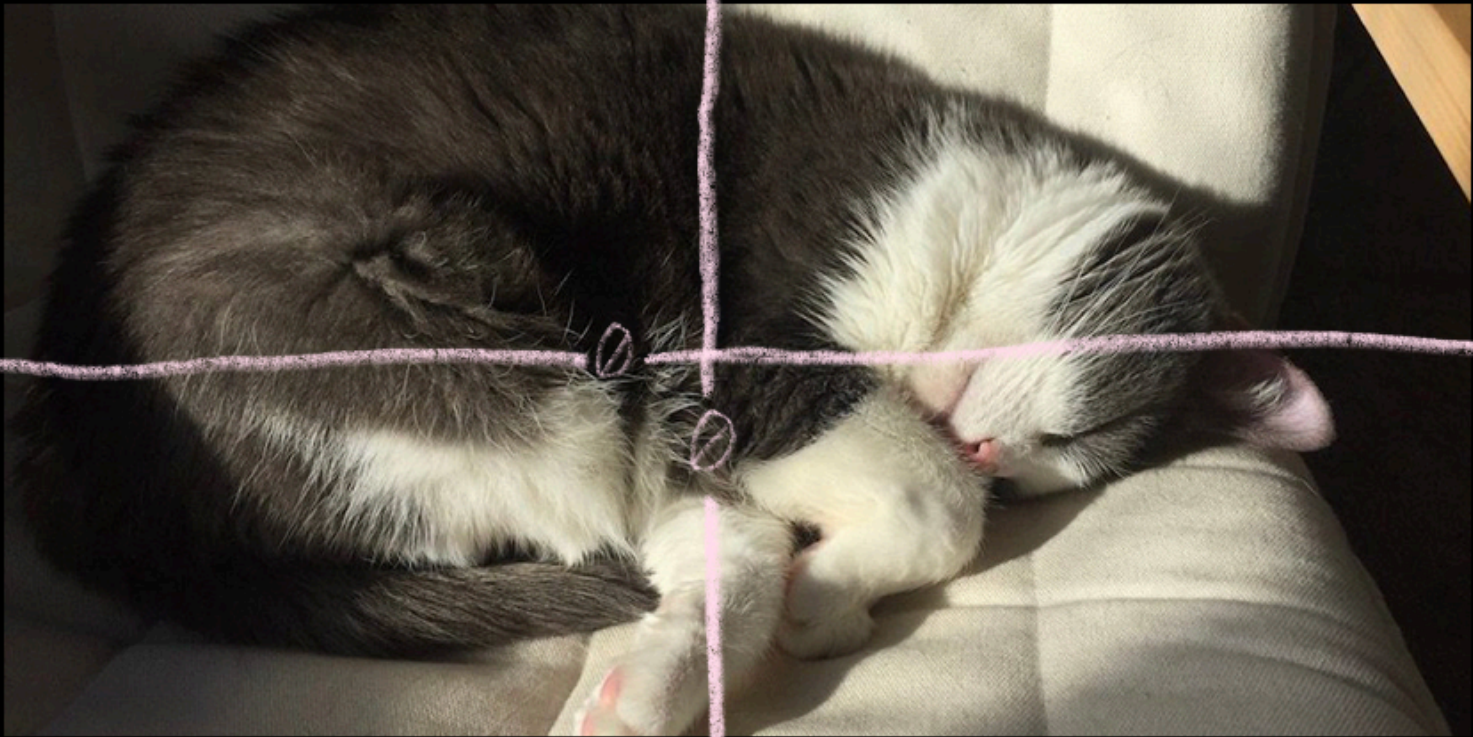
$\wedge$

$\text{very\_little\_grey}(\bar{x})$



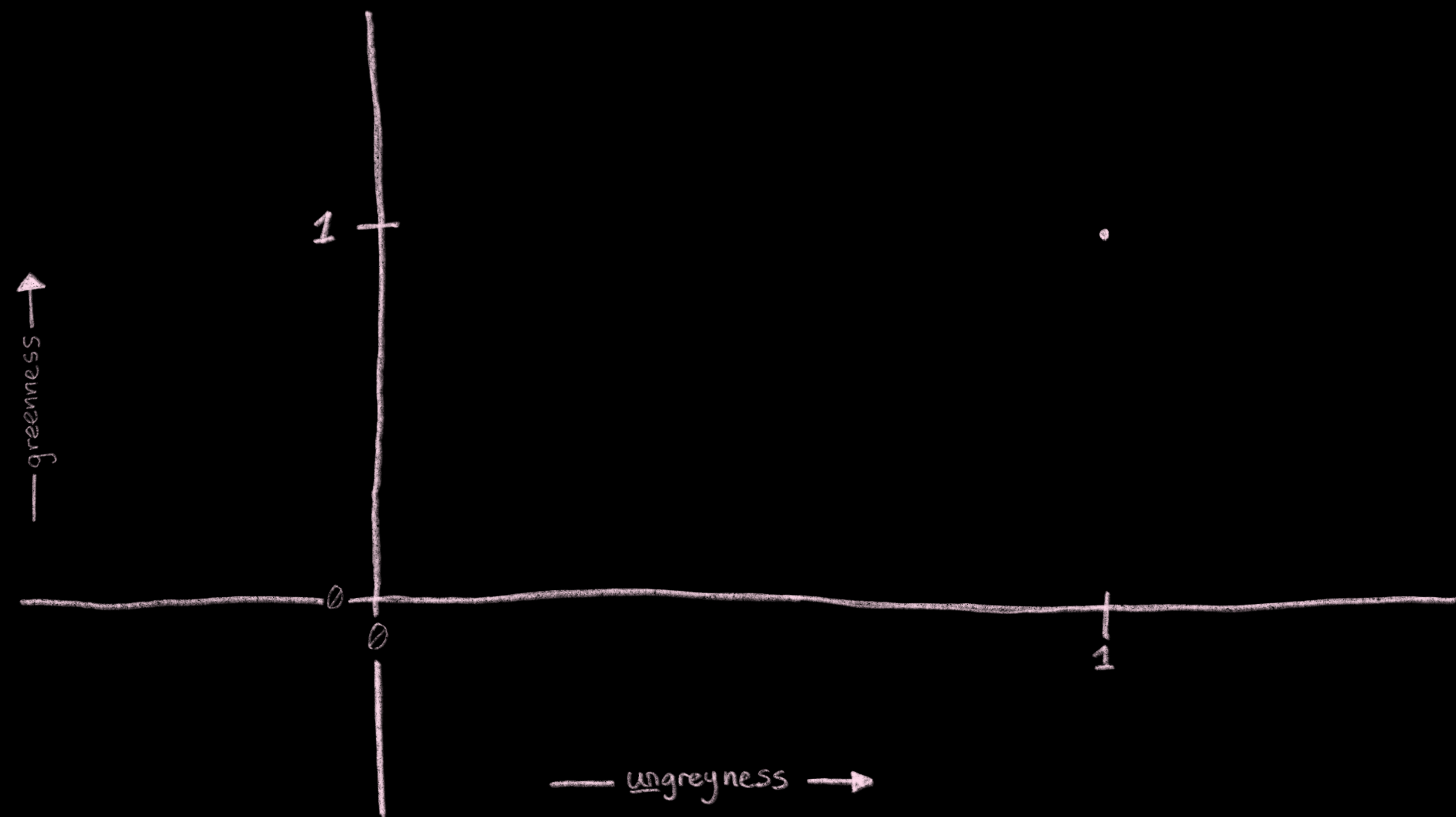


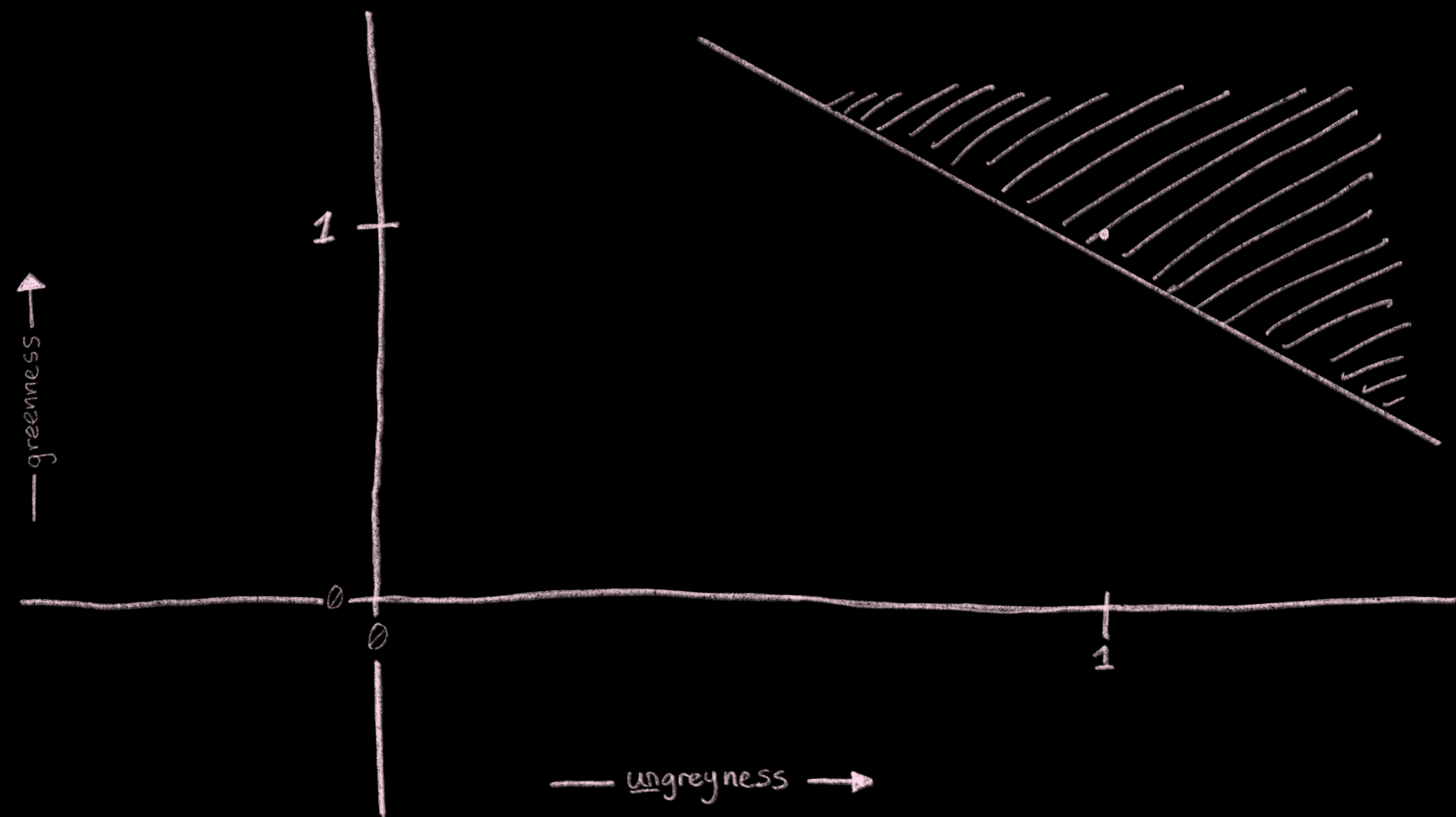
—greenness—→



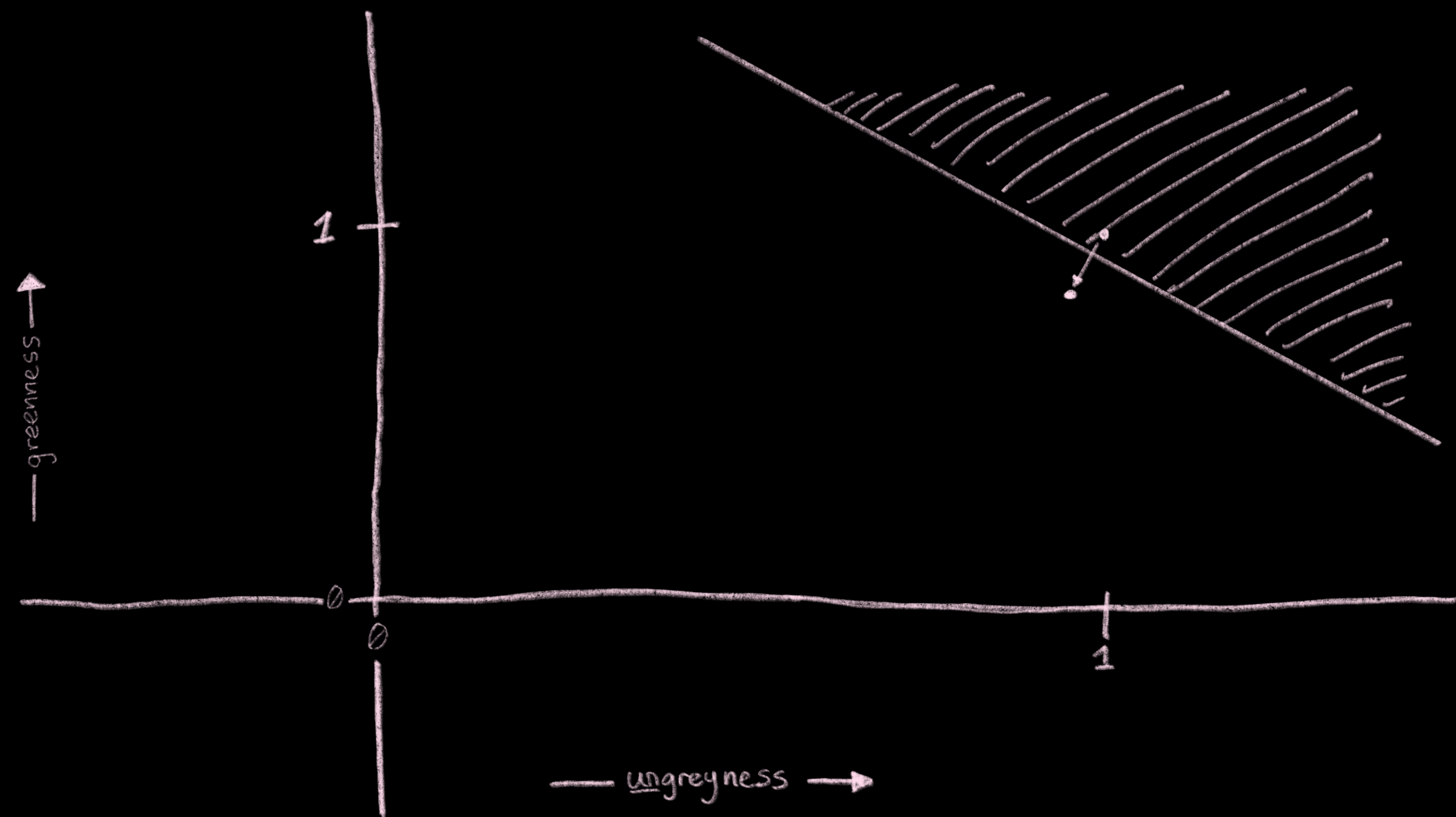
—ungrey<sup>n</sup>ness—→



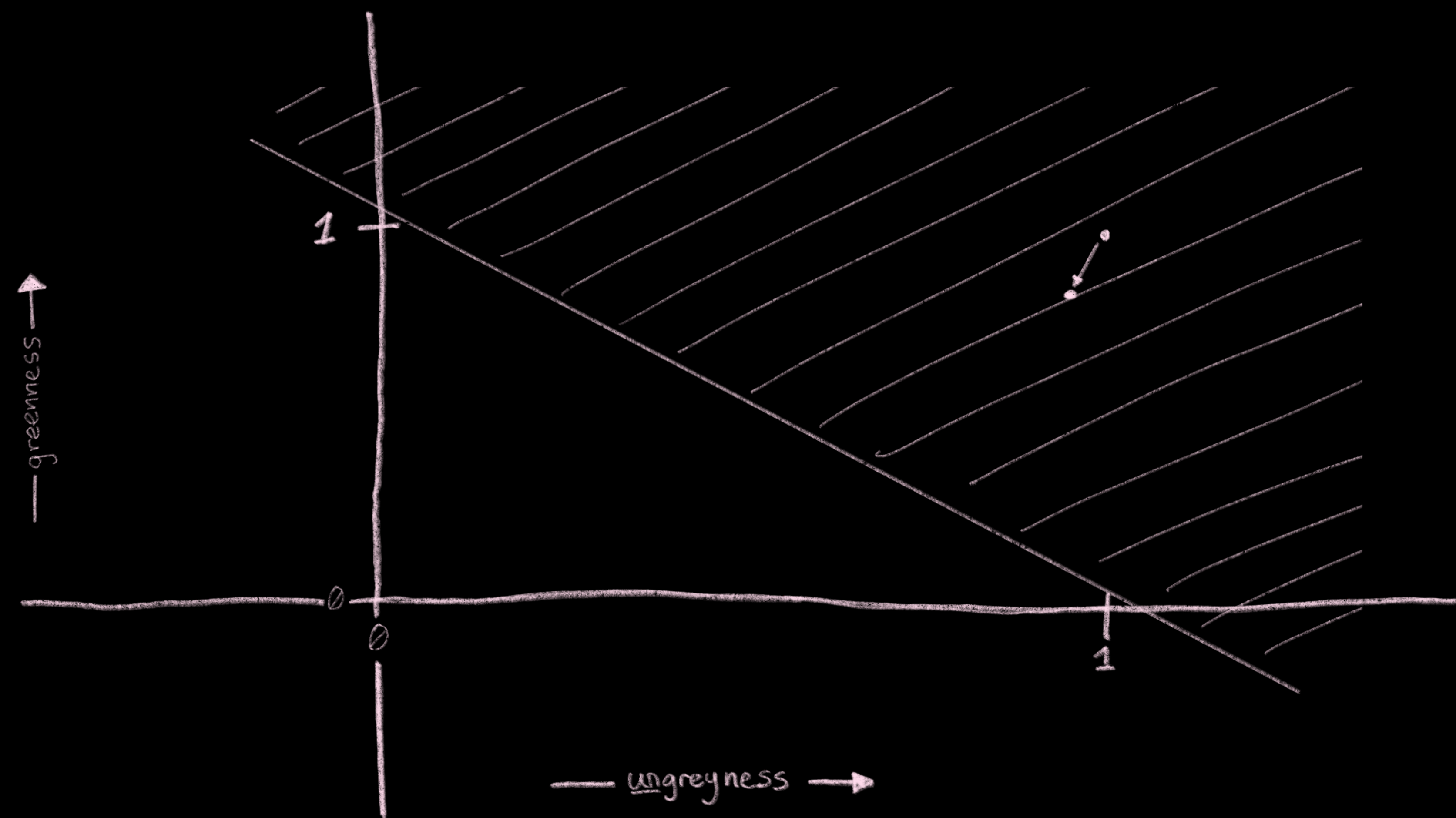




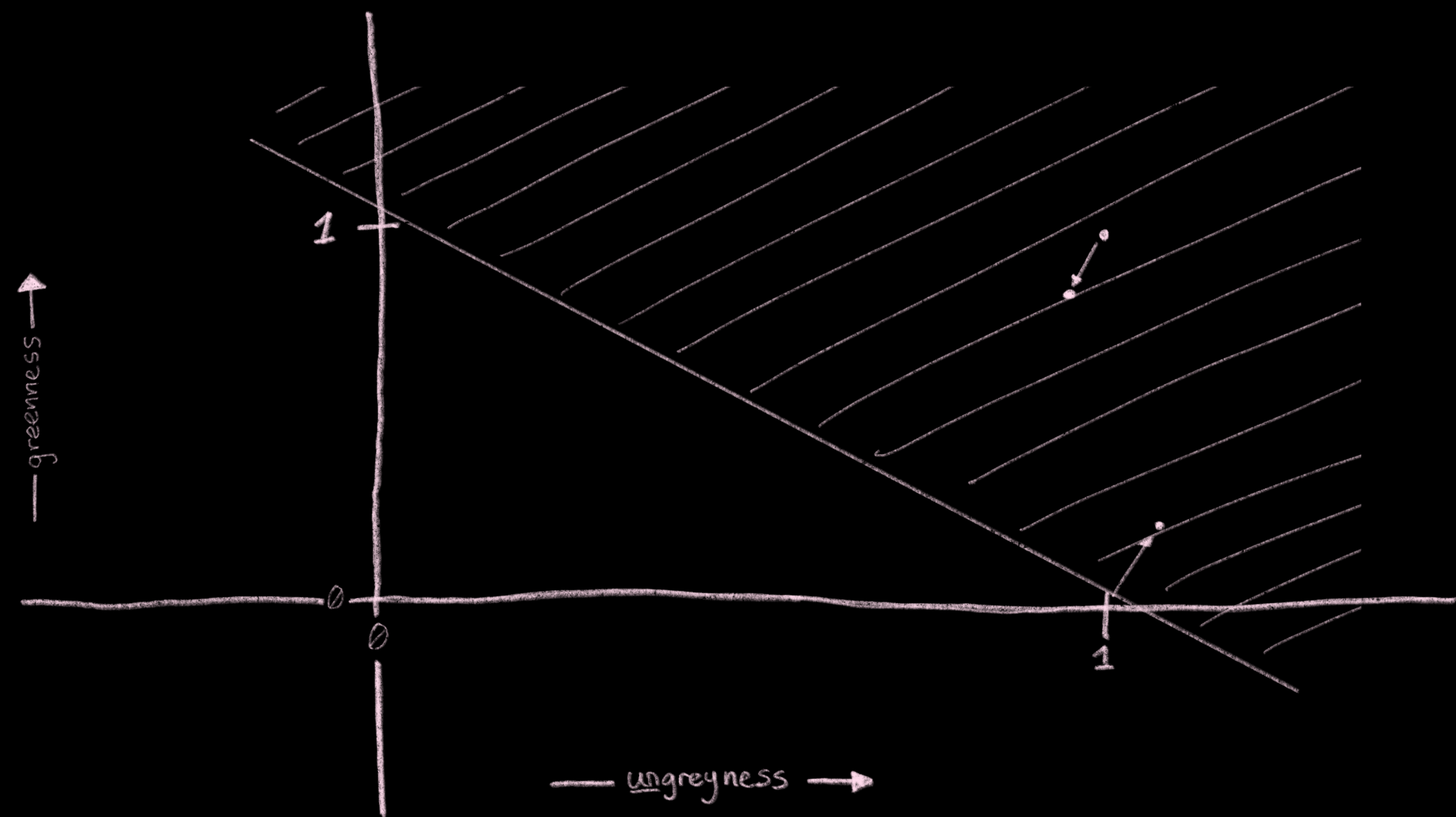










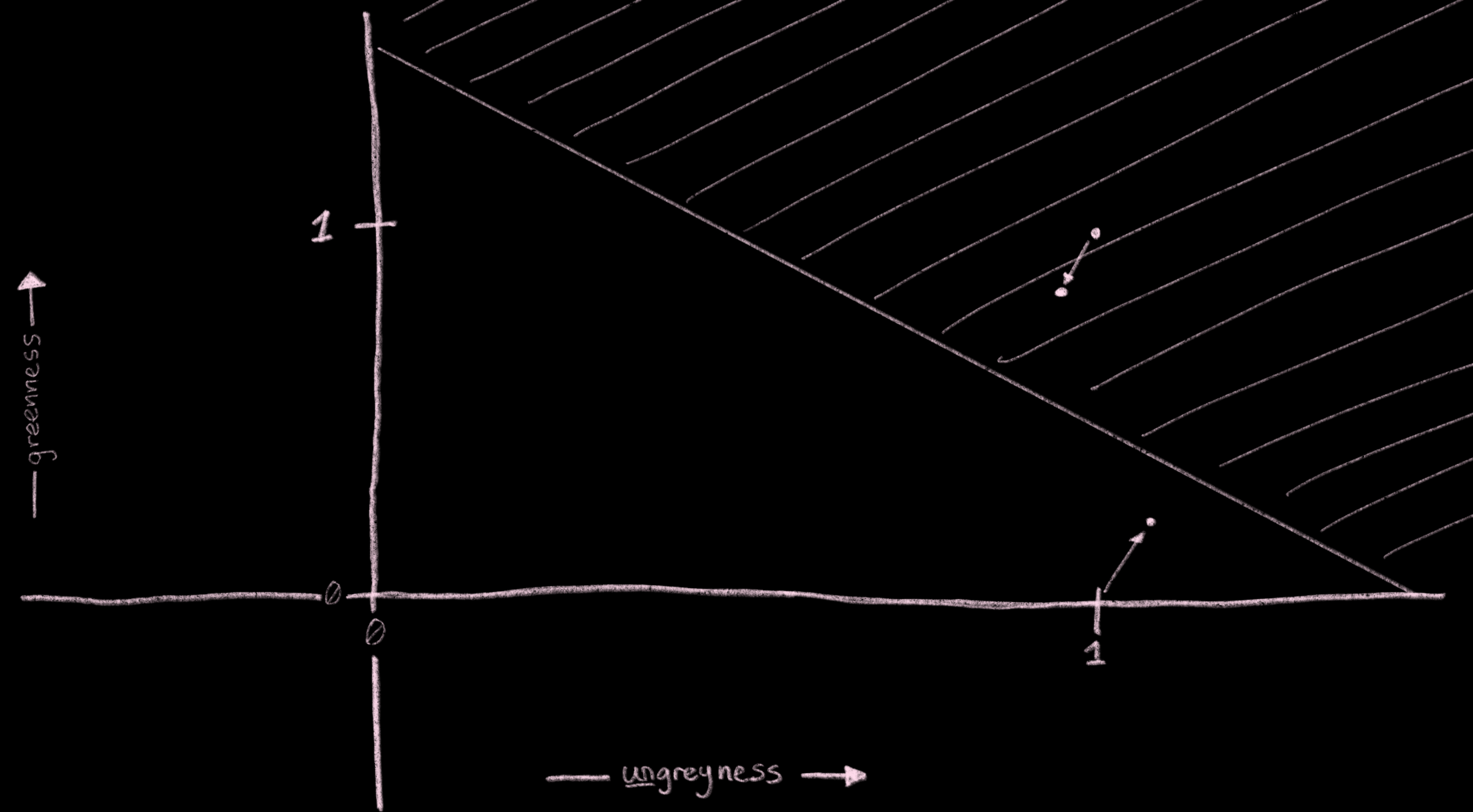






- DOG?







WHAT HAVE WE SEEN SO FAR?



WHAT'S STILL TO COME?



IS OUR  
DATA  
GOOD?

DID WE  
FIND THE  
RIGHT  
FEATURES?

ARE WE  
ROBUST  
AROUND  
THESE  
FEATURES?



IS OUR  
DATA  
GOOD?

DID WE  
FIND THE  
RIGHT  
FEATURES?

ARE WE  
ROBUST  
AROUND  
THESE  
FEATURES?



ROBUSTNESS AS A REFINEMENT TYPE



# WHAT'S A REFINEMENT TYPE?

A type refined with an SMT-checkable predicate.

```
let  $\mathbb{R}^+$  = (x: $\mathbb{R}$  {0.0R ≤ x}) // positive reals
```

```
let _ = 4.0R :  $\mathbb{R}^+$  // 0.0R ≤ 4.0R
```

```
let vector a n = (xs:list a {length xs = n}) // lists of length n
```

```
let _ = [0.5R; 1.0R] : vector  $\mathbb{R}^+$  2 // length [0.5R; 1.0R] = 2
```



**classify** :  $(x_1 \rightarrow \mathbb{R}) \rightarrow (x_2 \rightarrow \mathbb{R}) \rightarrow (y : \mathbb{R})$

**classify**  $x_1 \ x_2 = f (w_1 x_1 + w_2 x_2 - b)$



**classify** :  $(x_1 \rightarrow \mathbb{R}) \rightarrow (x_2 \rightarrow \mathbb{R}) \rightarrow (y : \mathbb{R})$

**classify**  $x_1 \ x_2 = S \ (0.5x_1 + 0.5x_2 - 0.9)$

$$S \ x = \begin{cases} 1, & \text{if } x \geq 0 \\ 0, & \text{otherwise} \end{cases}$$



```
val model : network (*with*) 2 (*inputs*) 1 (*output*) 1 (*layer*)
let model = NLast // makes single-layer network
  { weights      = [[0.5R]; [0.5R]]
  ; biases      = [-0.9R]
  ; activation = Threshold
  }

val classify : (x1 : ℝ) → (x2 : ℝ) → (y : ℝ)
let classify x1 x2 = run model [x1; x2]
```

```
let  $\varepsilon$  = 0.1R // how big are tiny steps?
```

```
val doggy : (x :  $\mathbb{R}$ )  $\rightarrow$  bool
```

```
let doggy x = 1.0R -  $\varepsilon$   $\leq$  x && x  $\leq$  1.0R +  $\varepsilon$ 
```

```
val _ = (x1 :  $\mathbb{R}\{\text{doggy } x1\}$ )
```

```
     $\rightarrow$  (x2 :  $\mathbb{R}\{\text{doggy } x2\}$ )
```

```
     $\rightarrow$  (y :  $\mathbb{R}\{y = 1.0R\}$ )
```

```
val _ = classify
```



```
(define-fun classify ((x1 Real) (x2 Real)) Real
  (ite (>= (- (+ (* x1 0.5) (* x2 0.5)) 0.9) 0.0) 1.0 0.0))
(define-fun doggy ((x Real)) Bool (and (<= 0.9 x) (<= x 1.1)))
(assert (forall ((x1 Real) (x2 Real))
  (=> (and (doggy x1) (doggy x2)) (= (classify x1 x2) 1.0))))
(check-sat)
```

> sat ;; it works! your network is totally robust! gj!

SO IT WORKS.

BUT DOES IT WORK?



①

SOLVERS  
DON'T DO  
NON-LINEAR  
ARITHMETIC

②

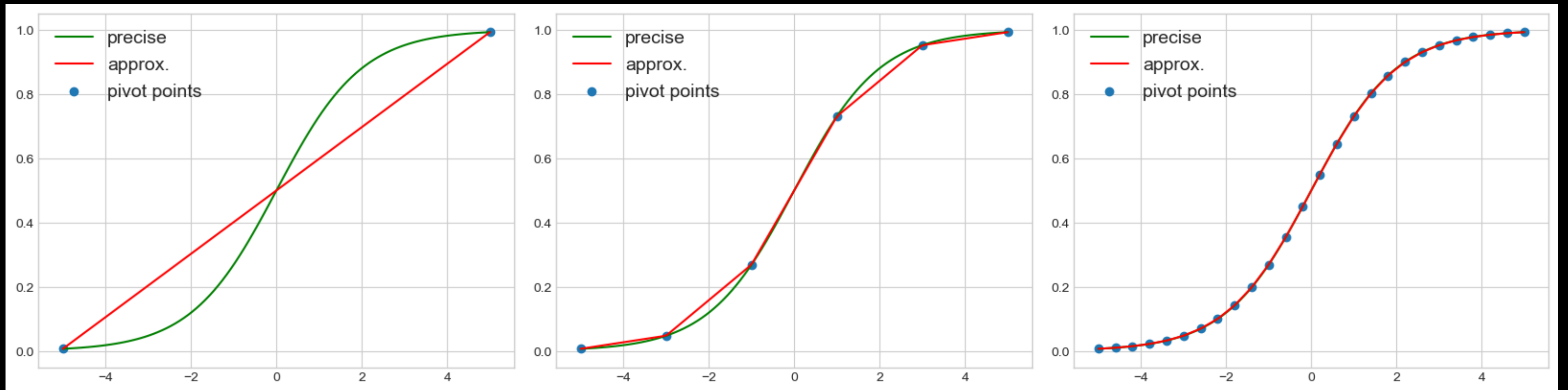
INTEGRATION  
WITH  $F^*$   
INTRODUCES A  
SIGNIFICANT  
SLOWDOWN

③

SOLVERS  
DON'T SCALE  
TO REALISTIC  
SIZES

# ① SOLVERS DON'T DO NON-LINEAR ARITHMETIC

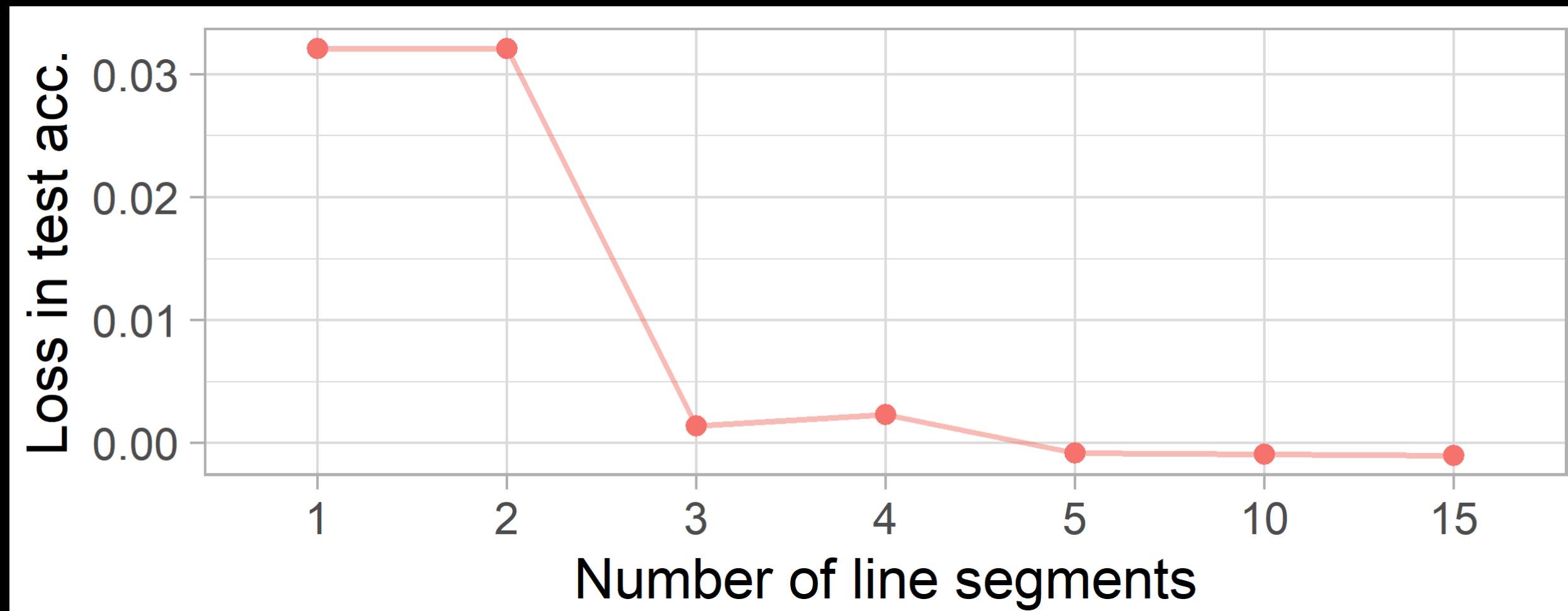
Let's make our activation functions linear!



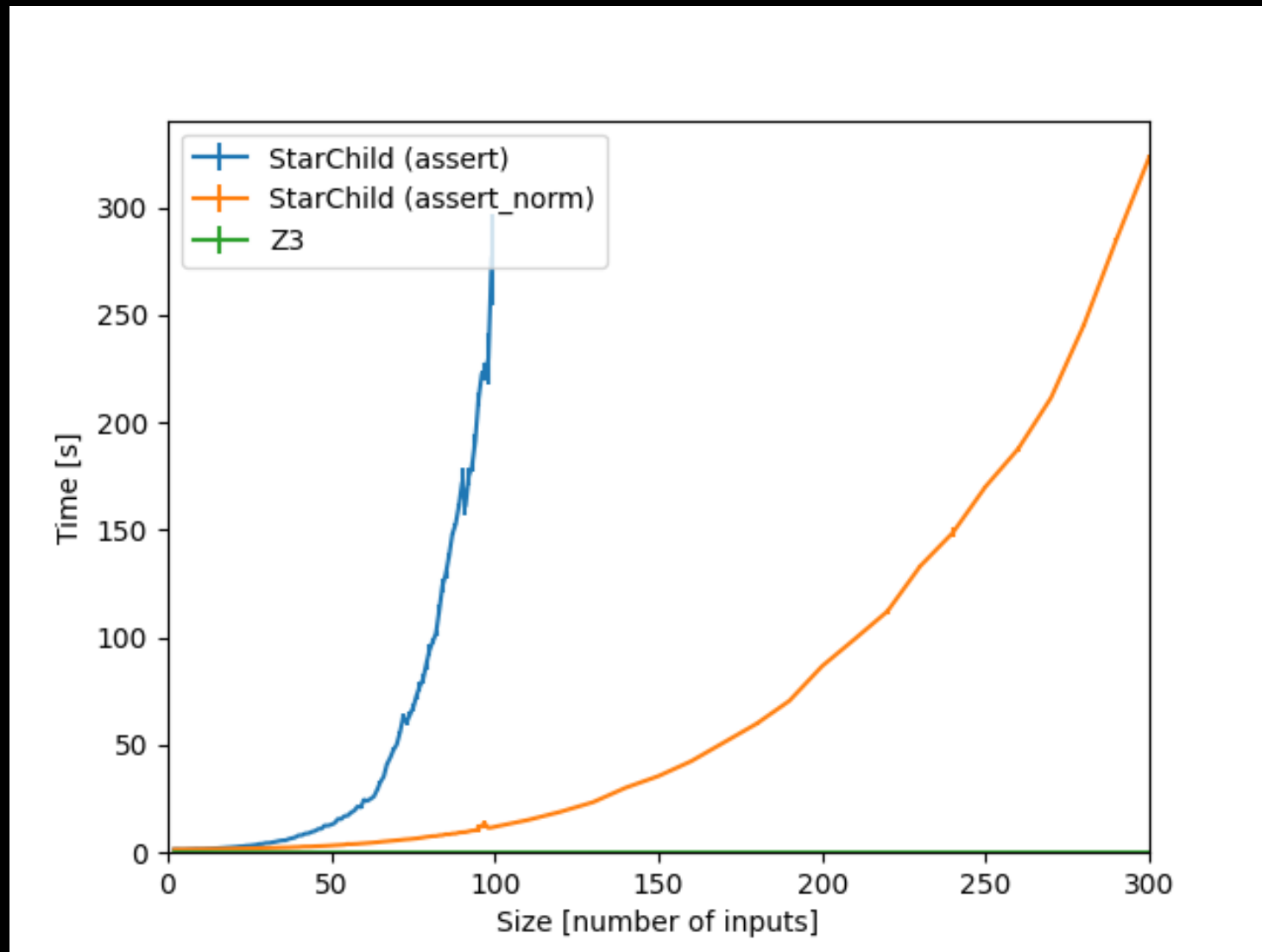


# ① SOLVERS DON'T DO NON-LINEAR ARITHMETIC

Train with tanh, run with linear approximation!



## ② INTEGRATION WITH $F^*$ INTRODUCES A SIGNIFICANT SLOWDOWN



Ahh! An exponential!

Don't make Z3 do  
reduction!

Don't tell Z3 about  
data-types.

(Unless you have to.)



# ③ SOLVERS DON'T SCALE TO REALISTIC SIZES

Z3 ignores tons of structure!

MetiTarski solves exponentials!

nnenum solves ReLUs!

Marabou solves piecewise-linear functions!

# ROBUSTNESS AS A REFINEMENT TYPE

- encode robustness as a refinement type
- leverage existing integration with solvers
- lightweight verification of robustness

but

- need to improve integration with solvers
- need more flexibility in choosing solvers