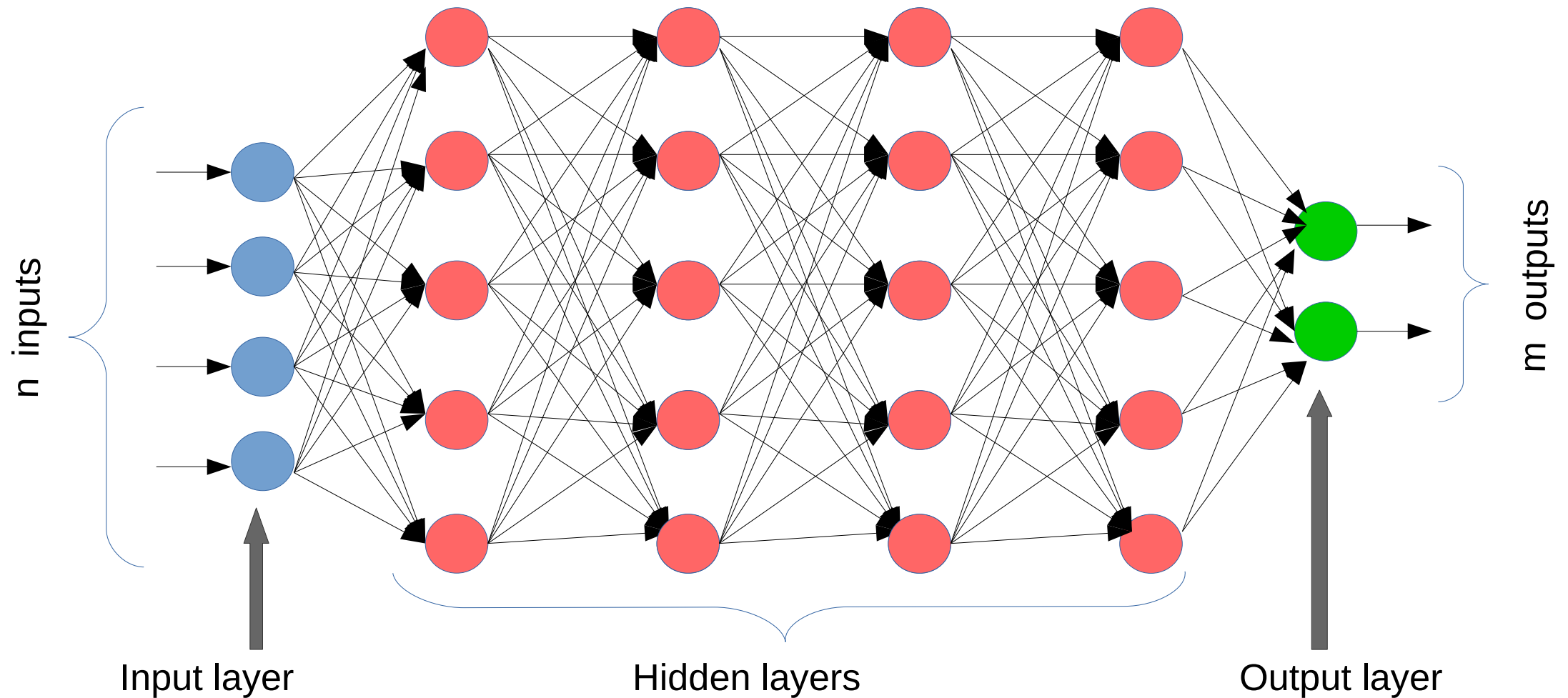


# **CS620: FM in ML**

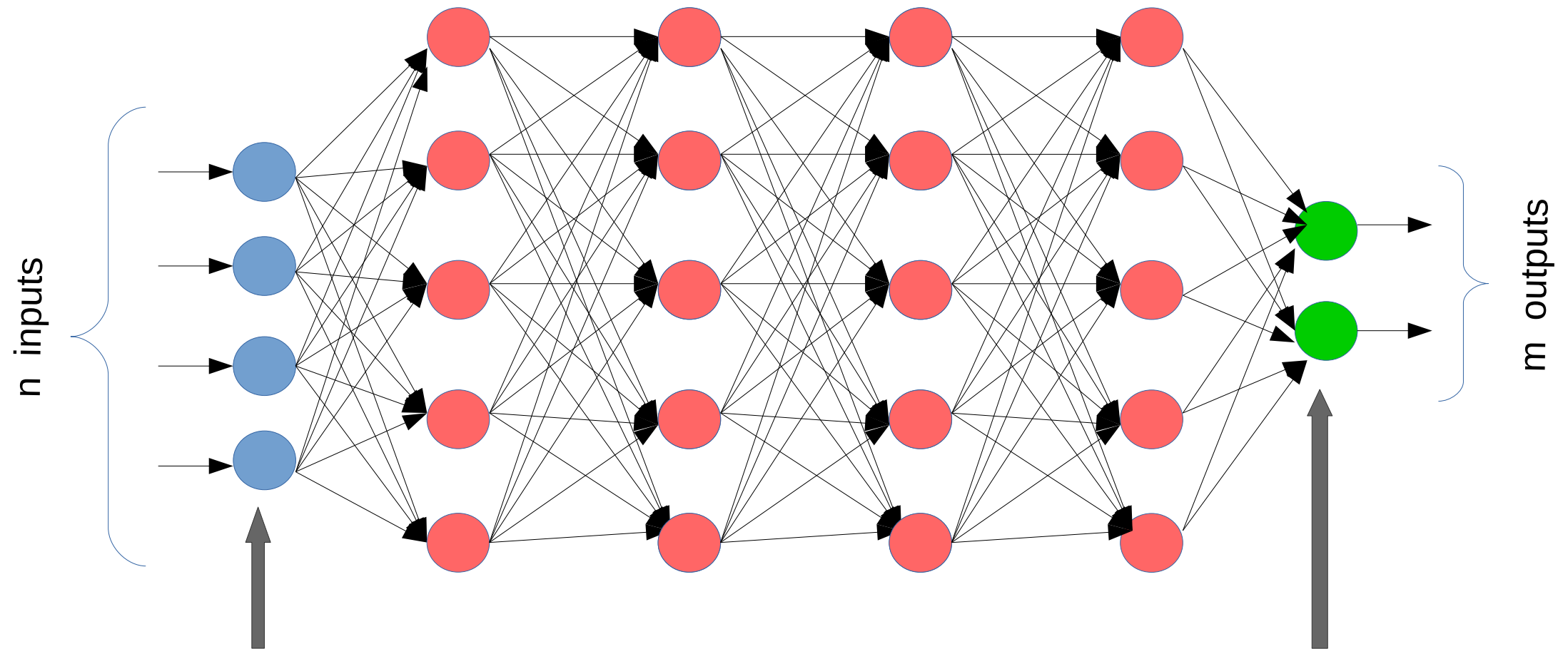
## **Specifying Properties of Neural Networks (Week 3)**

Supratik Chakraborty

# A Typical Neural Network



# A Typical Neural Network



Input layer

Output layer

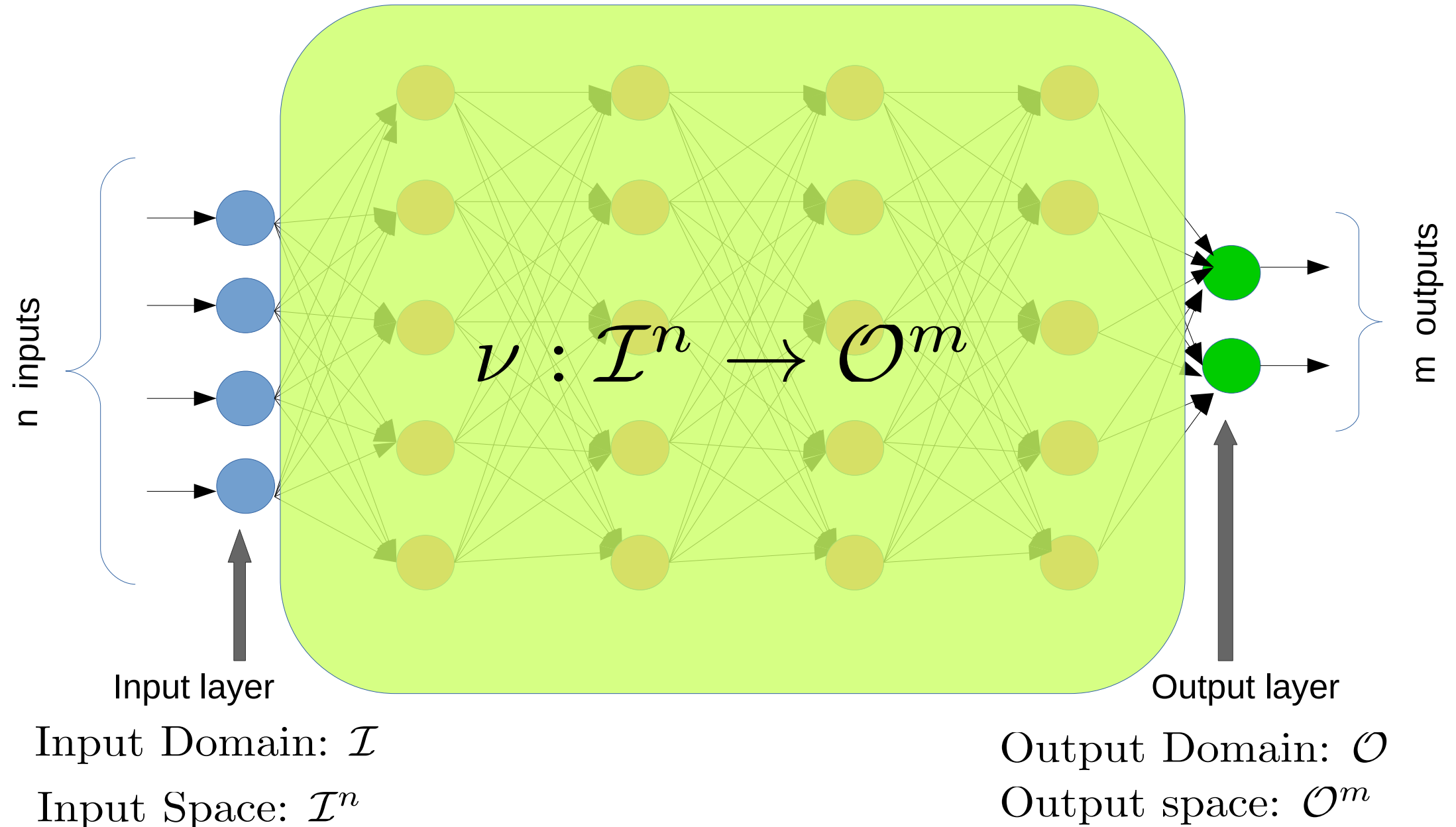
Input Domain:  $\mathcal{I}$

Input Space:  $\mathcal{I}^n$

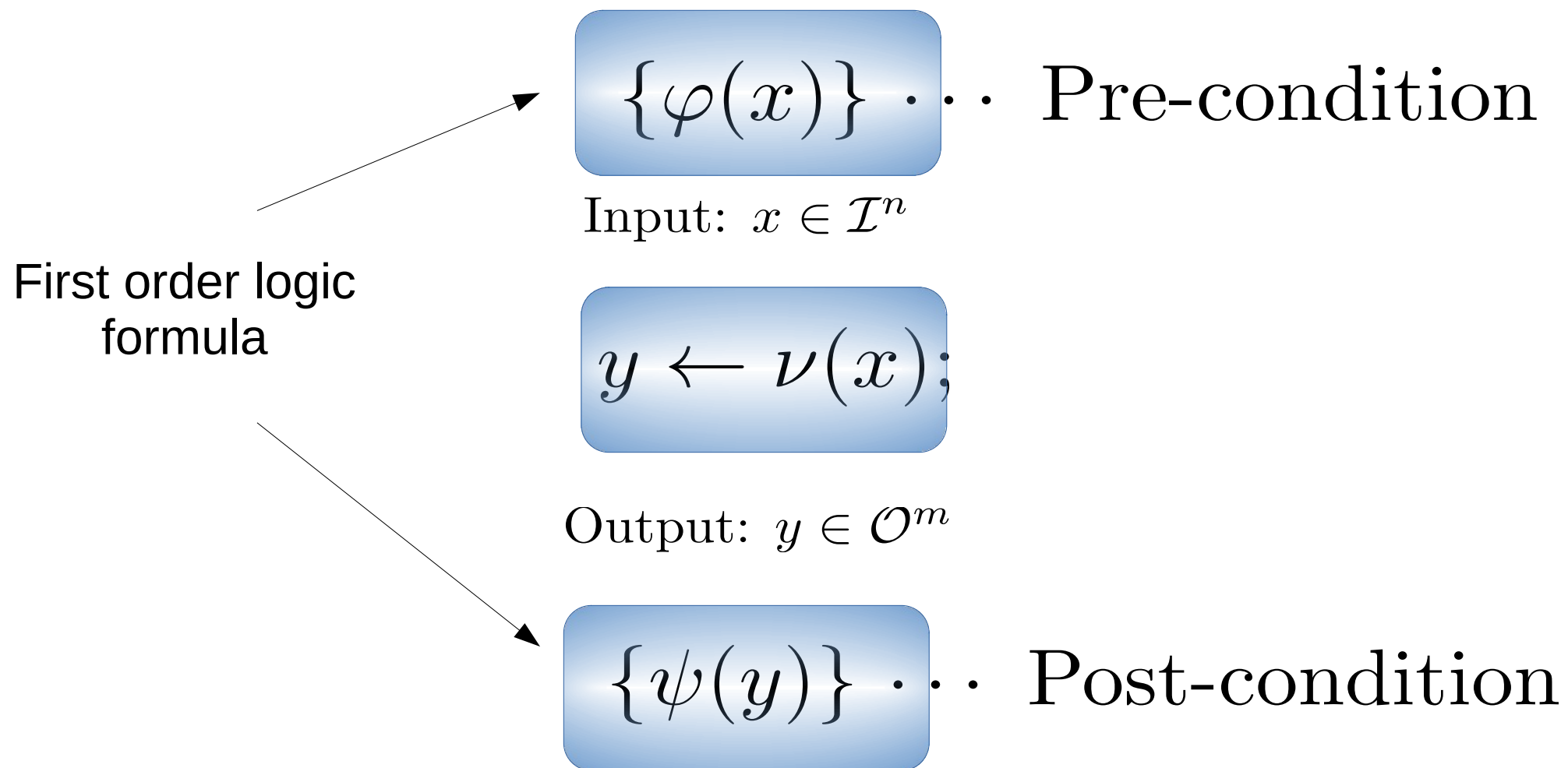
Output Domain:  $\mathcal{O}$

Output space:  $\mathcal{O}^m$

# A Typical Neural Network



# A Transformative Program



Hoare triples similar to those used in program verification

# Semantics of Hoare Triple

$\{\varphi(x)\} \dots$  Pre-condition

$y \leftarrow \nu(x); \dots$  "Program"

$\{\psi(y)\} \dots$  Post-condition

Validity of Hoare triple

If  $x$  satisfies  $\varphi(x)$ ,  
"program" terminates and encounters no memory exception,  
then output  $y$  always satisfies  $\psi(y)$

# Property Specification Example 1



"panda"

57.7% confidence

+  $\epsilon$



=



"gibbon"

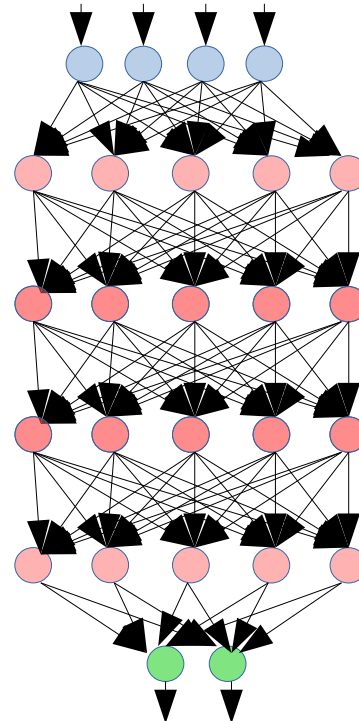
99.3% confidence

Source: Goodfellow, Shlens, Szegedy, "Explaining and Harnessing Adversarial Examples", 2015

Wish to specify that the above never happens  
for a given image, for a specified max perturbation

# Property Specification Example 1

Specified image:  $x^*$



Score for panda:  $p$

Score for something else:  $g$

$$\{\|x - x^*\| \leq \varepsilon\}$$

Max perturbation of input

$$(p, g) \leftarrow \nu(x)$$

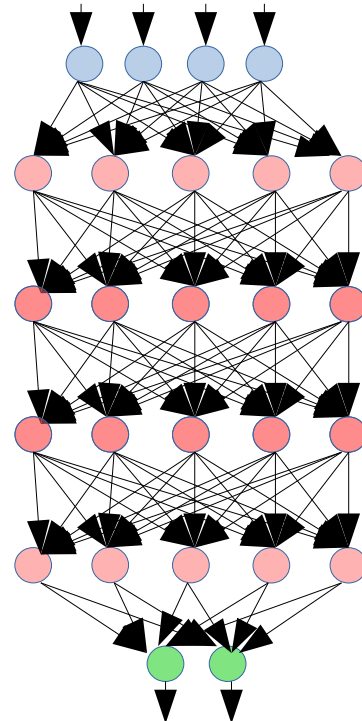
Separation threshold for  
“confident” classification

$$\{p > g + \delta\}$$



# Property Specification Example 1

Specified image:  $x^*$



Score for panda:  $p$

Score for something else:  $g$

$$\{\|x - x^*\| \leq \varepsilon\}$$

$$\bigwedge_{i=1}^N (|r_i - r_i^*| \leq \varepsilon_r) \wedge \\ \bigwedge_{i=1}^N (|g_i - g_i^*| \leq \varepsilon_g) \wedge \\ \bigwedge_{i=1}^N (|b_i - b_i^*| \leq \varepsilon_b)$$

$$(p, g) \leftarrow \nu(x);$$

$$\{p > g + \delta\}$$

# Spec as a logical requirement

$$\forall r_1 \forall g_1 \forall b_1 \cdots \forall r_N \forall g_N \forall b_N \forall p \forall g$$

$$\left( \begin{array}{l} \bigwedge_{i=1}^N (|r_i - r_i^*| \leq \varepsilon_r) \wedge \\ \bigwedge_{i=1}^N (|g_i - g_i^*| \leq \varepsilon_g) \wedge \\ \bigwedge_{i=1}^N (|b_i - b_i^*| \leq \varepsilon_b) \wedge \\ (p, g) = \nu(r_1, g_1, b_1, \dots, r_N, g_N, b_N) \\ \\ \implies \\ p > g + \delta \end{array} \right)$$

$$\{\|x - x^*\| \leq \varepsilon\}$$

$$\bigwedge_{i=1}^N (|r_i - r_i^*| \leq \varepsilon_r) \wedge \\ \bigwedge_{i=1}^N (|g_i - g_i^*| \leq \varepsilon_g) \wedge \\ \bigwedge_{i=1}^N (|b_i - b_i^*| \leq \varepsilon_b)$$

$$(p, g) \leftarrow \nu(x);$$

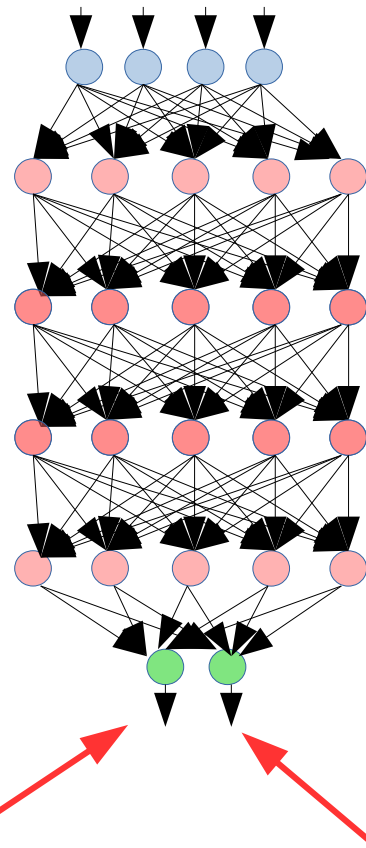
$$\{p > g + \delta\}$$

**A logical implication**

# Property Specification Example 2

Given two arbitrary images that differ within prescribed limits, the network must never “confidently” classify them differently

Arbitrary image  $x$



Score for class 1:  $s_1$

Score for class 2:  $s_2$

$$\{\|x - x^*\| \leq \varepsilon\}$$

$$(s_1, s_2) \leftarrow \nu(x);$$
$$(s_1^*, s_2^*) \leftarrow \nu(x^*);$$

$$\left\{ \begin{array}{l} (s_1 > s_2 + \delta) \implies (s_1^* > s_2^* + \delta) \wedge \\ (s_2 > s_1 + \delta) \implies (s_2^* > s_1^* + \delta) \end{array} \right\}$$

# Property Specification Example 2



## Pause n Reflect

Given two images that differ within prescribed limits, the network must never “confidently” classify them differently

$$\{\|x - x^*\| \leq \varepsilon\}$$

$$\begin{aligned}(s_1, s_2) &\leftarrow \nu(x); \\ (s_1^*, s_2^*) &\leftarrow \nu(x^*);\end{aligned}$$

Are there any unintended consequences of the specification?

Can a neural network satisfying the specification do anything meaningful?

How easy/hard is it to design a neural network satisfying this specification?

$$\left\{ \begin{aligned} (s_1 > s_2 + \delta) &\implies (s_1^* > s_2^* + \delta) \wedge \\ (s_2 > s_1 + \delta) &\implies (s_2^* > s_1^* + \delta) \end{aligned} \right\}$$

# Spec as a logical requirement

$$\forall r_1 \dots \forall b_N \forall r_1^* \dots \forall b_N^* \forall s_1 \forall s_2 \forall s_1^* \forall s_2^*$$

$$\bigwedge_{i=1}^N (|r_i - r_i^*| \leq \varepsilon_r) \wedge$$

$$\bigwedge_{i=1}^N (|g_i - g_i^*| \leq \varepsilon_g) \wedge$$

$$\bigwedge_{i=1}^N (|b_i - b_i^*| \leq \varepsilon_b) \wedge$$

$$(s_1, s_2) = \nu(r_1, g_1, b_1, \dots, r_N, g_N, b_N) \wedge$$

$$(s_1^*, s_2^*) = \nu(r_1^*, g_1^*, b_1^*, \dots, r_N^*, g_N^*, b_N^*) \wedge$$

$$\implies$$

$$(s_1 > s_2 + \delta) \implies (s_1^* > s_2^* + \delta) \wedge$$

$$(s_2 > s_1 + \delta) \implies (s_2^* > s_1^* + \delta)$$

$$\{\|x - x^*\| \leq \varepsilon\}$$

$$\bigwedge_{i=1}^N (|r_i - r_i^*| \leq \varepsilon_r) \wedge$$

$$\bigwedge_{i=1}^N (|g_i - g_i^*| \leq \varepsilon_g) \wedge$$

$$\bigwedge_{i=1}^N (|b_i - b_i^*| \leq \varepsilon_b)$$

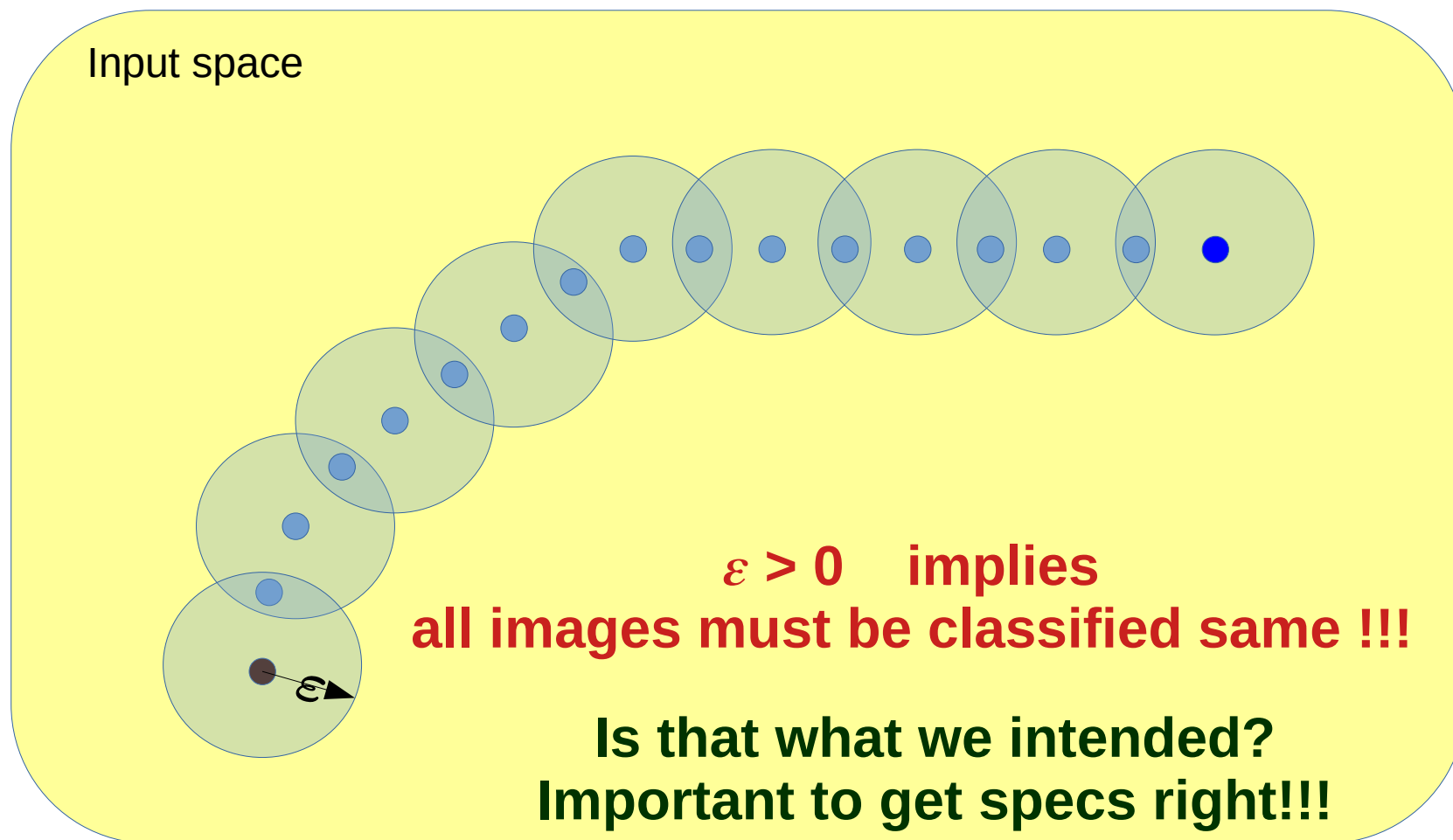
$$(s_1, s_2) \leftarrow \nu(x);$$

$$(s_1^*, s_2^*) \leftarrow \nu(x^*);$$

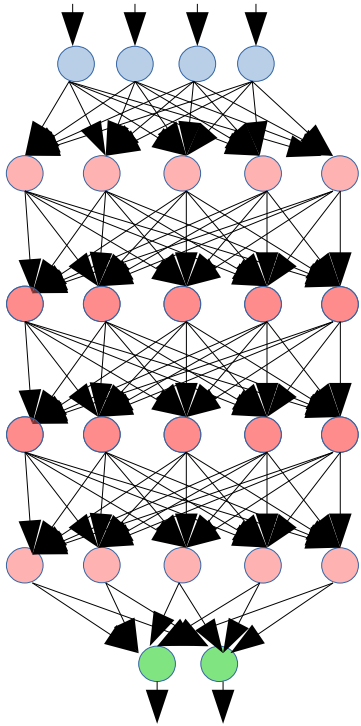
$$\left\{ \right.$$

# Problem with Specification 2

Pick any two arbitrary images in the input space



# Taking a step back to re-look



Arbitrary input

$$\{\|x - x^*\| \leq \varepsilon\}$$

Specific input

$$(p, g) \leftarrow \nu(x)$$

$$\{p > g + \delta\}$$

Arbitrary input

$$\{\|x - x^*\| \leq \varepsilon\}$$

Arbitrary input

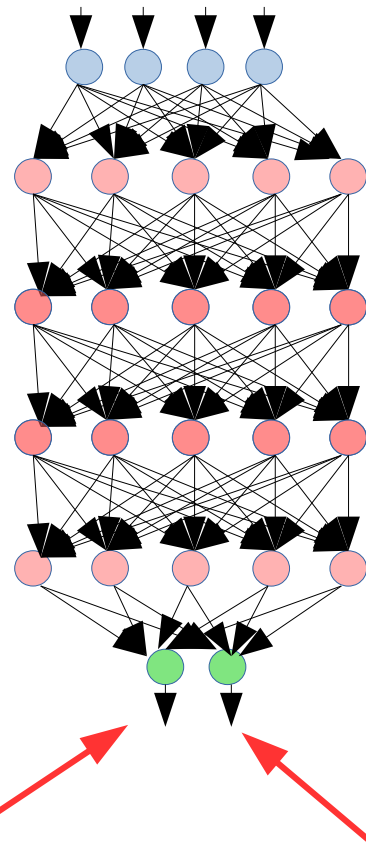
$$\begin{aligned} (s_1, s_2) &\leftarrow \nu(x); \\ (s_1^*, s_2^*) &\leftarrow \nu(x^*); \end{aligned}$$

$$\left\{ \begin{aligned} (s_1 > s_2 + \delta) &\implies (s_1^* > s_2^* + \delta) \wedge \\ (s_2 > s_1 + \delta) &\implies (s_2^* > s_1^* + \delta) \end{aligned} \right\}$$

# Attempting a Fix

Given two arbitrary images that differ within prescribed limits, the network must never “confidently” classify them differently

Arbitrary image  $x$



Score for class 1:  $s_1$

Score for class 2:  $s_2$

$$\{\|x - x^*\| \leq \varepsilon\}$$

$$(s_1, s_2) \leftarrow \nu(x);$$
$$(s_1^*, s_2^*) \leftarrow \nu(x^*);$$

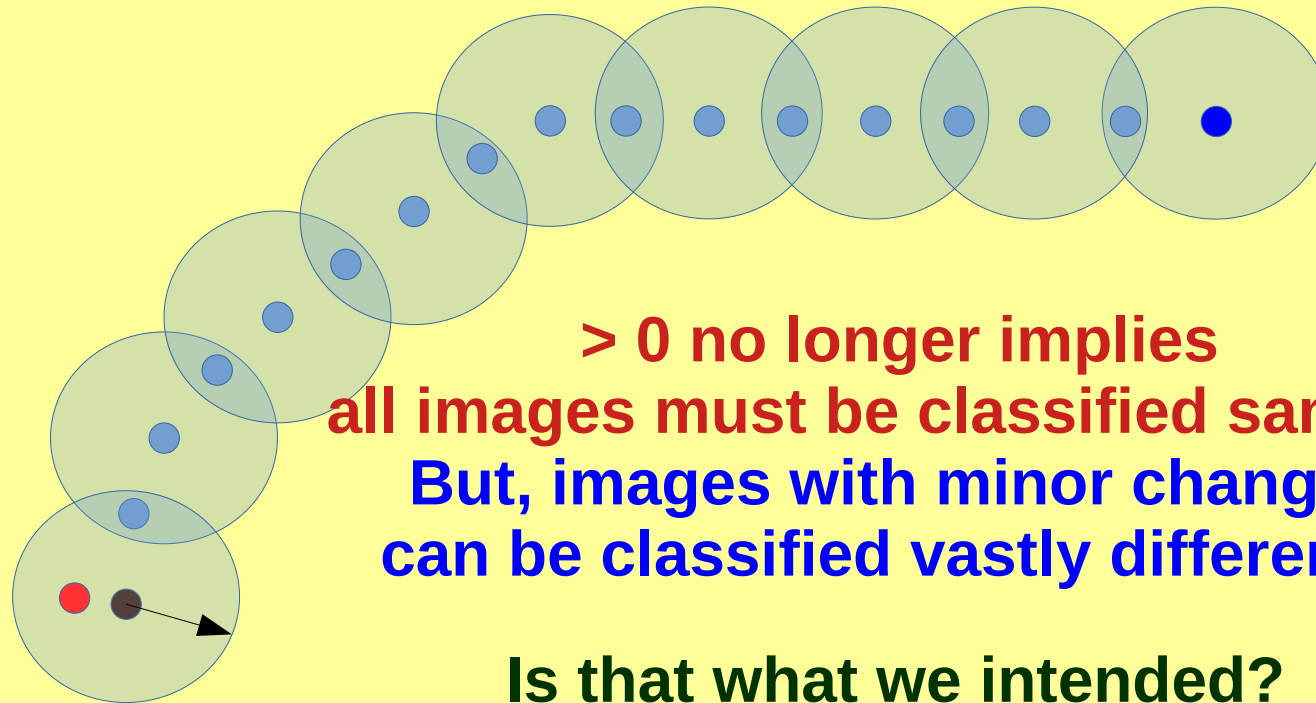
$$\left\{ \begin{array}{l} (s_1 > s_2 + \delta) \implies (s_2^* \leq s_1^* + \delta) \wedge \\ (s_2 > s_1 + \delta) \implies (s_1^* \leq s_2^* + \delta) \end{array} \right\}$$



# Did It Fix?

Pick any two arbitrary images in the input space

Input space



**> 0 no longer implies  
all images must be classified same !!!**

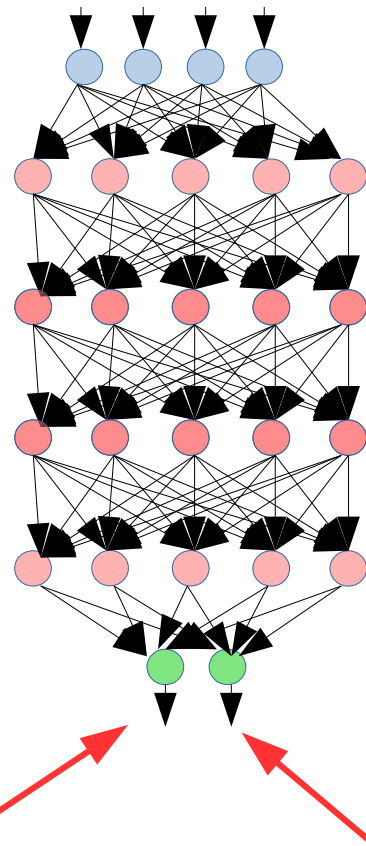
**But, images with minor changes  
can be classified vastly differently**

**Is that what we intended?  
Important to get specs right!!!**

# Property Specification Example 2

## Second attempt!

Given two arbitrary images that differ pixel-wise within prescribed limits and have “similar” semantic features, the network must never “confidently” classify them differently



$$\left\{ \left( \|x - x^*\| \leq \varepsilon \right) \wedge \left( \sigma(x) \approx \sigma(x^*) \right) \right\}$$

$$\begin{aligned} (s_1, s_2) &\leftarrow \nu(x); \\ (s_1^*, s_2^*) &\leftarrow \nu(x^*); \end{aligned}$$

$$\left\{ \begin{aligned} (s_1 > s_2 + \delta) &\implies (s_1^* > s_2^* + \delta) \wedge \\ (s_2 > s_1 + \delta) &\implies (s_2^* > s_1^* + \delta) \end{aligned} \right\}$$

Score for class 1:  $s_1$

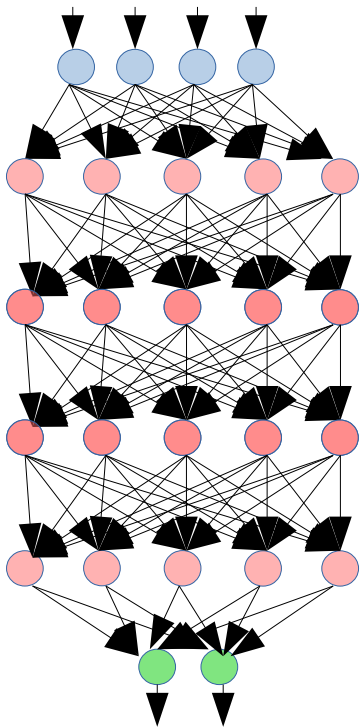
Score for class 2:  $s_2$

# Property Specification Example 2

## Second attempt!

Given two arbitrary images that differ pixel-wise within prescribed limits and have “similar” semantic features, the network must never “confidently” classify them differently

$$\left\{ \left( \|x - x^*\| \leq \varepsilon \right) \wedge \left( \sigma(x) \approx \sigma(x^*) \right) \right\}$$



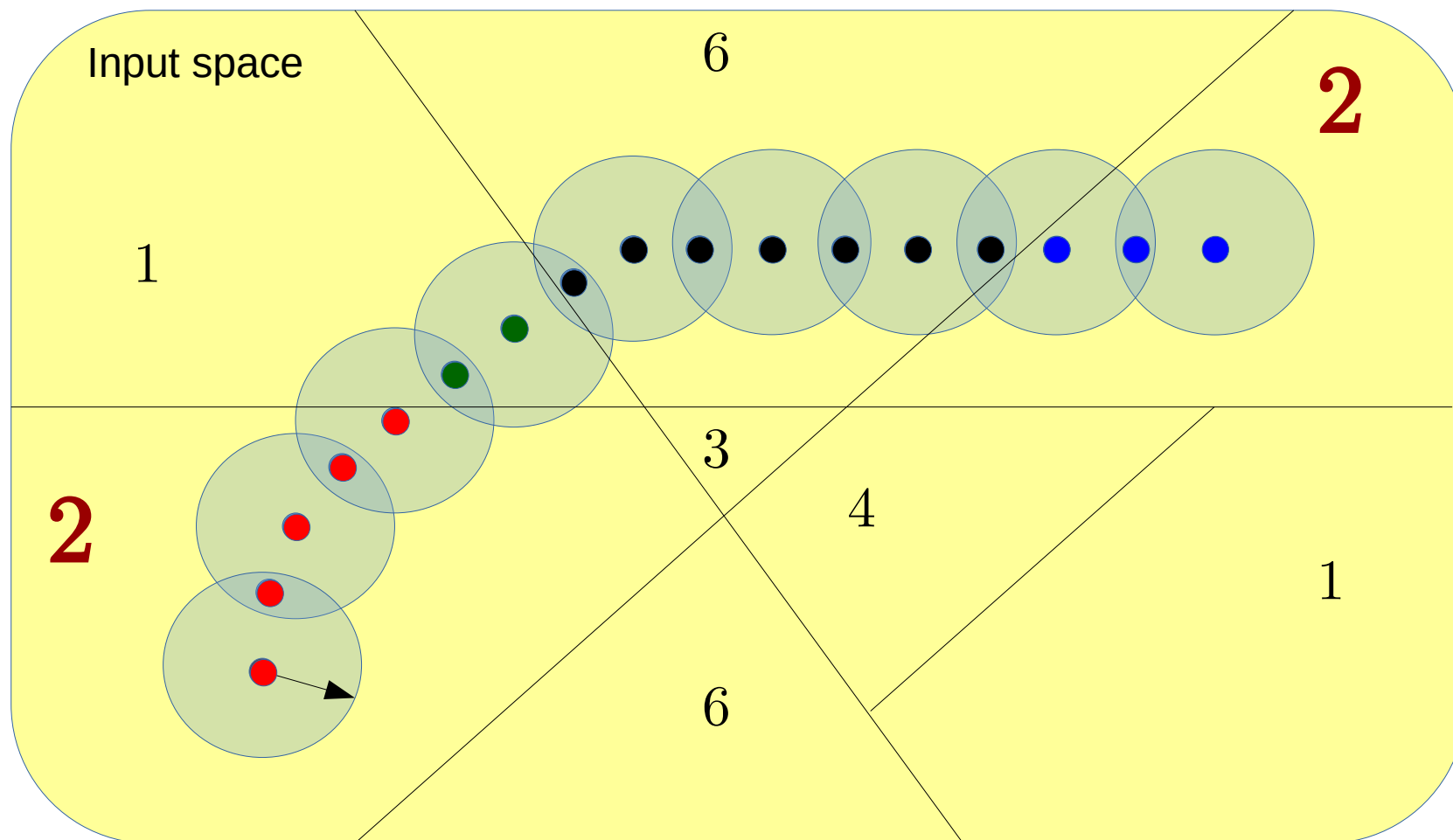
User-defined semantic features,  
Not necessarily network-defined

$$\begin{aligned} (s_1, s_2) &\leftarrow \nu(x); \\ (s_1^*, s_2^*) &\leftarrow \nu(x^*); \end{aligned}$$

$$\left\{ \begin{aligned} (s_1 > s_2 + \delta) &\implies (s_1^* > s_2^* + \delta) \wedge \\ (s_2 > s_1 + \delta) &\implies (s_2^* > s_1^* + \delta) \end{aligned} \right\}$$

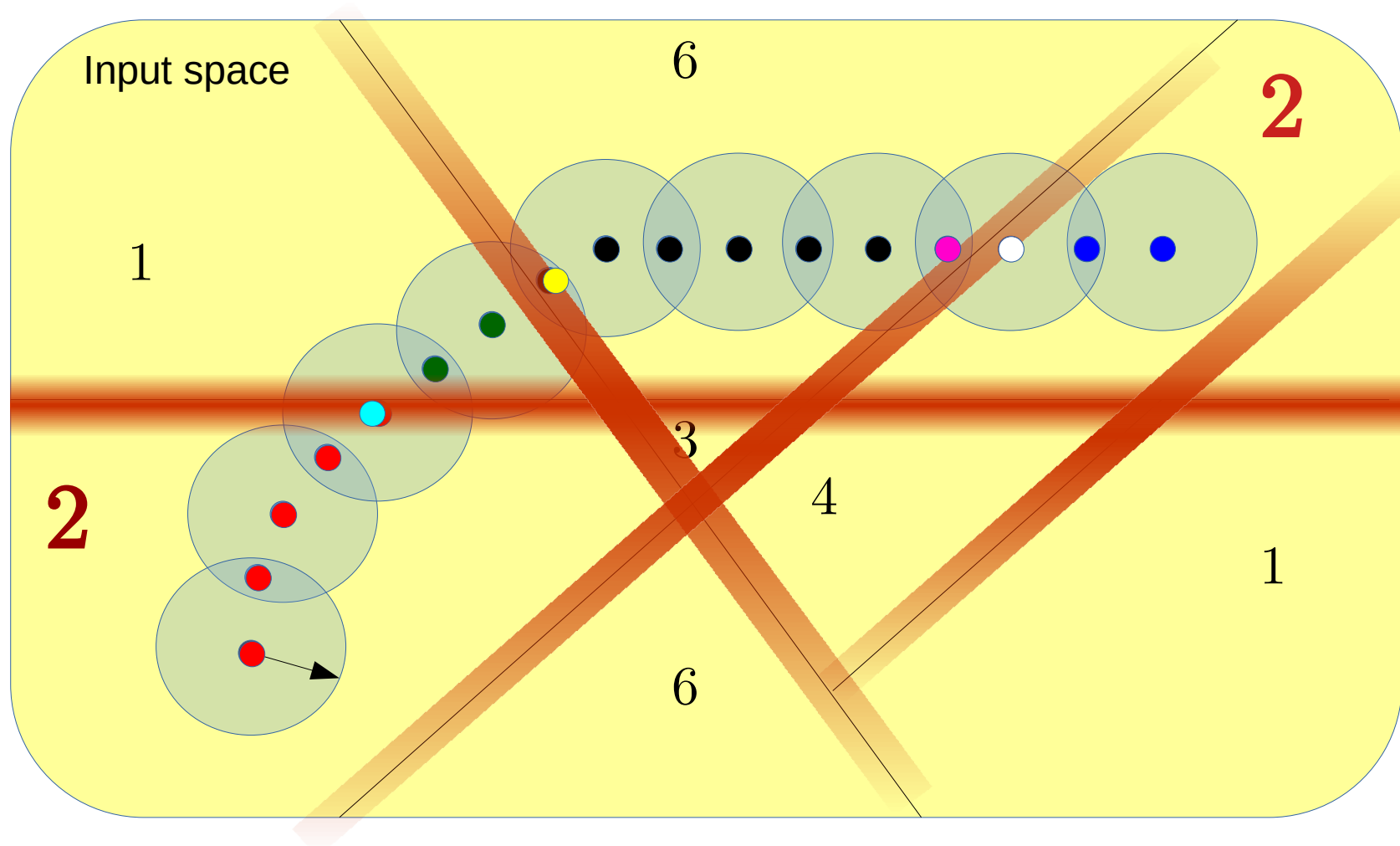
# Possibilities with New Spec

Pick any two arbitrary images in the input space



# Possibilities with Newer Spec

Pick any two arbitrary images in the input space



# Property Specification Example 2

## Third attempt!

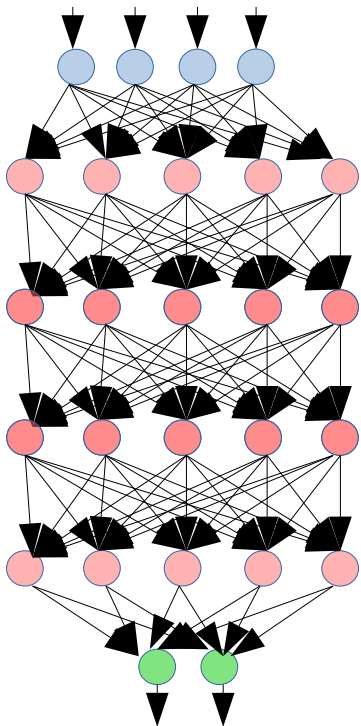
Given two arbitrary images that differ pixel-wise within prescribed limits and have “similar” semantic features, the network must produce “similar” classifications

$$\left\{ \left( \|x - x^*\| \leq \varepsilon \right) \wedge \left( \sigma(x) \approx \sigma(x^*) \right) \right\}$$

$$\begin{aligned} (s_1, s_2) &\leftarrow \nu(x); \\ (s_1^*, s_2^*) &\leftarrow \nu(x^*); \end{aligned}$$

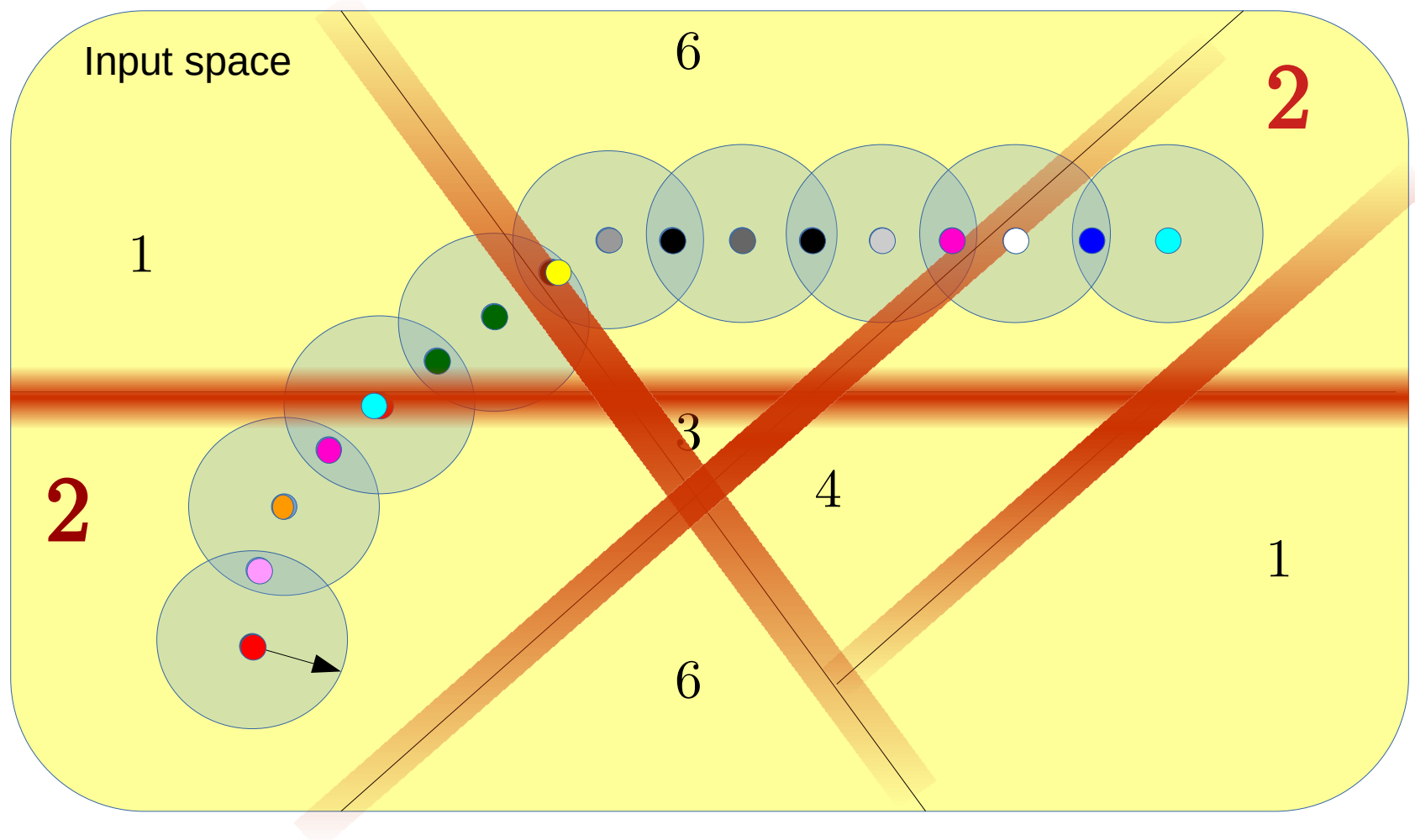
$$\left\{ \lambda(s_1, s_2) \simeq \lambda(s_1^*, s_2^*) \right\}$$

Network-defined  
labeling function:  
“final” layer(s)



# Possibilities with New Spec

Pick any two arbitrary images in the input space



# Property Specification



**Pause n Reflect**

**Why is it so hard to get specifications right?**

**Is it easier to arrive at**

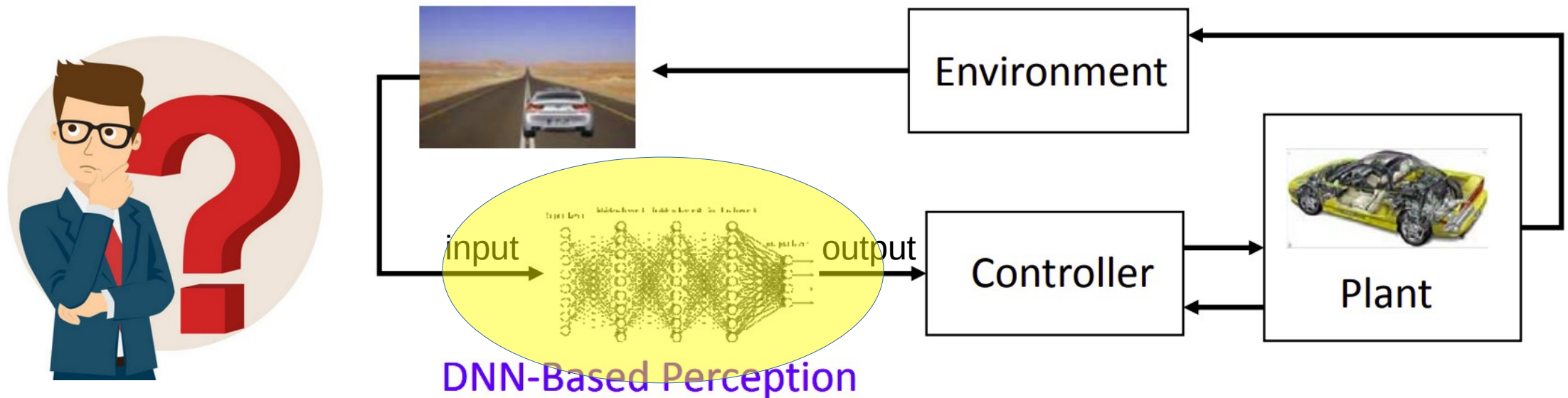
**THE RIGHT SPECIFICATION that covers all aspects of behaviour**

**OR**

**A bunch of sub-specifications that cover parts of the behaviour space?**



# A Day In The Life of A “Specifier”



Source: Seshia et al, Formal Verification of Deep Neural Networks, 2018

Collect a bunch of **desired/undesired** (input, output) pairs

- Not necessarily what DNN is actually doing
- Instead, what DNN's environment “expects” it to do



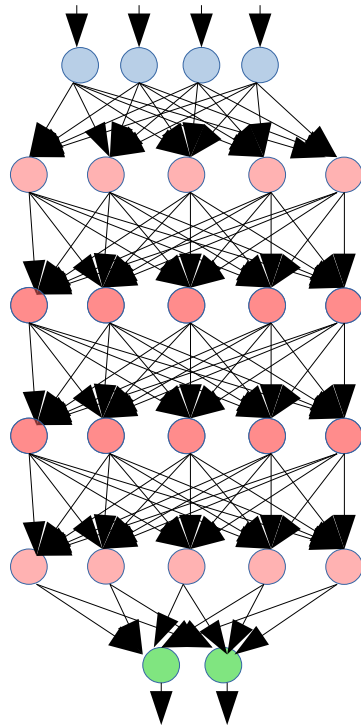
Is there a **formalizable relation** between inputs and desired outputs?

- Did we miss out corner cases?
- Sufficiently constrained to preclude all undesired behaviour?
- Sufficiently relaxed to allow all desired behaviour?

# Input-Output Relation: How hard is it to formalize?

## Self-driving car

Image (road scene)

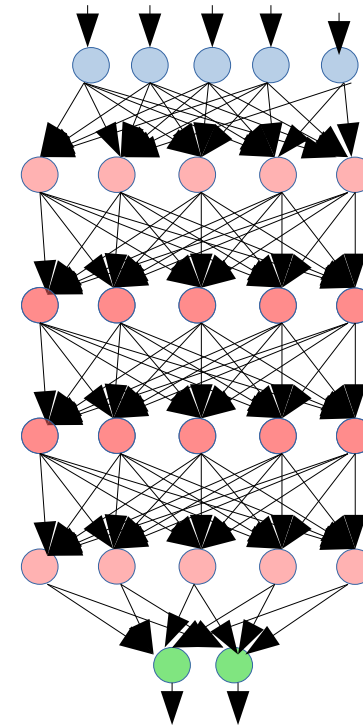


Perceptual Spec

“Too congested to accelerate”

## Unmanned drone

Flight parameters

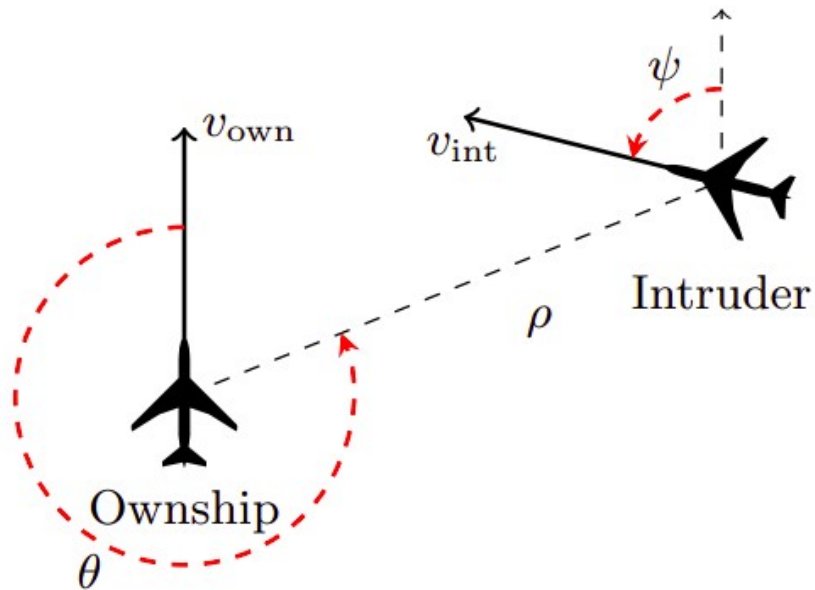


Non-Perceptual  
Spec

Score  
(Horizontal Advisory)

# Non-Perceptual DNN Specs

## ACAS-Xu



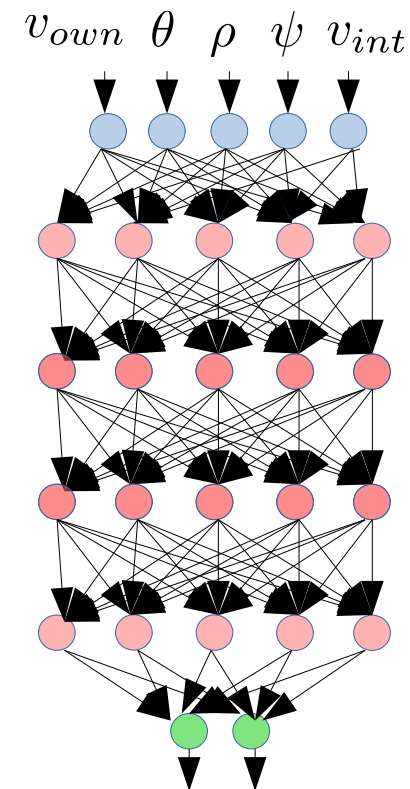
$$\{(\rho \geq 55947.691ft) \wedge (v_{own} \geq 1145ft/s) \wedge (v_{int} \leq 60ft/s)\}$$

$$\text{Score} \leftarrow \nu(\rho, v_{own}, v_{int}, \theta, \psi)$$

$$\{\text{Score}[\text{COC}] \leq 1500\}$$

Clear-of-Conflict

### Flight parameters



Score  
(Horizontal Advisory)

# Non-Perceptual DNN Specs

## ACAS-Xu

Rules for ACAS-Xu when directly implemented  
takes  $> 2\text{GB}$  memory

Flight parameters

45 Non-perceptual DNNs for same take  $< 3\text{MB}$  of memory

Having a good spec for a non-perceptual DNN  
doesn't make the DNN irrelevant !!!

Specs NOT SAME AS Rules

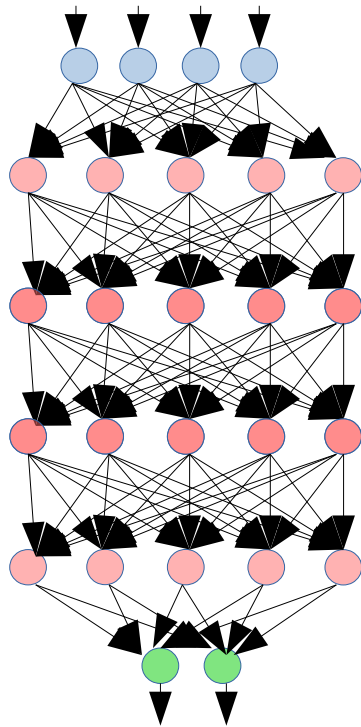
$\{\text{Score}[\text{COC}] \leq 1500\}$

Score  
(Horizontal Advisory)

# Perceptual DNN Specs

$$(r_1, g_1, b_1, \dots, r_N, g_N, b_N)$$

Image (road scene)



“Too congested to accelerate”

Good spec:

$$\text{CR}(r_1, g_1, b_1, \dots, r_N, g_N, b_N)$$

~~$\{ (r_1, g_1, b_1, \dots, r_N, g_N, b_N) : \text{image of congested road} \}$~~

$$y \leftarrow \nu(r_1, g_1, b_1, \dots, r_N, g_N, b_N);$$

$\{ y = \text{“Too congested to accelerate”} \}$

**$\text{CR}( \dots ) = \text{true iff road is “too congested to accelerate”}$**



**If we know CR (...), why design and train a DNN ???**

# Perceptual DNN Specs

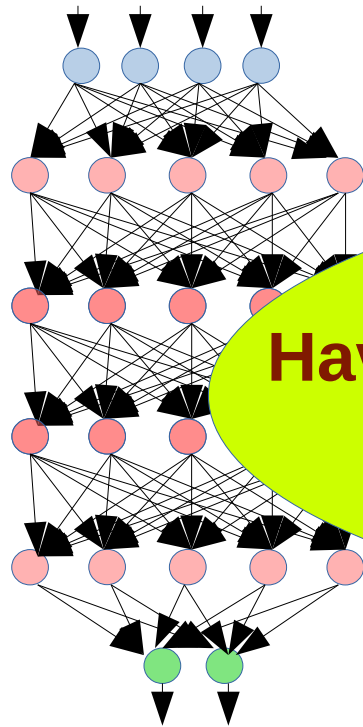
$$(r_1, g_1, b_1, \dots, r_N, g_N, b_N)$$

Image (road scene)

Good spec:

$$\text{CR}(r_1, g_1, b_1, \dots, r_N, g_N, b_N)$$

~~$\{ (r_1, g_1, b_1, \dots, r_N, g_N, b_N) : \text{image of congested road} \}$~~



**Having the ideal spec for a perceptual DNN  
would make the DNN irrelevant !!!**

$a_N, b_N);$

” }

sted to accelerate”

“Too congested to accelerate”



**If we know CR (...), why design and train a DNN ???**

# Specifying Properties of Perceptual DNNs



**Pause n Reflect**

**Are we in a chicken-and-egg conundrum for perceptual DNNs?**

**Is there any meaningful way out?**

**We can talk about robustness of classification w.r.t. a specific image**

**Can we specify anything formally beyond this?**

# Points to Ponder

**Are we in a chicken-and-egg conundrum for perceptual DNNs?**

**Is there any meaningful way out?**

**We can talk about robustness of classification w.r.t. a specific image**

**Can we specify anything formally beyond this?**

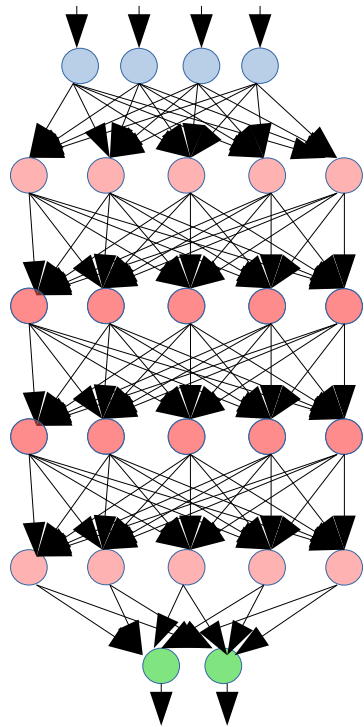
**Is it better to write a single all-encompassing spec or multiple sub-specs for different behavioural requirements?**



# Any Hope for Perceptual DNNs?

$$(r_1, g_1, b_1, \dots, r_N, g_N, b_N)$$

Image (road scene)

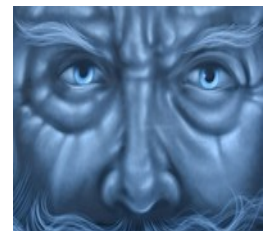
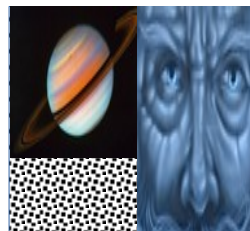
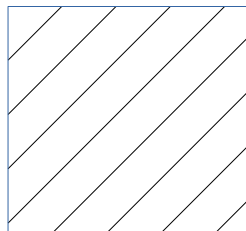


Input is

**High dimensional, large input space**

| Input Space |

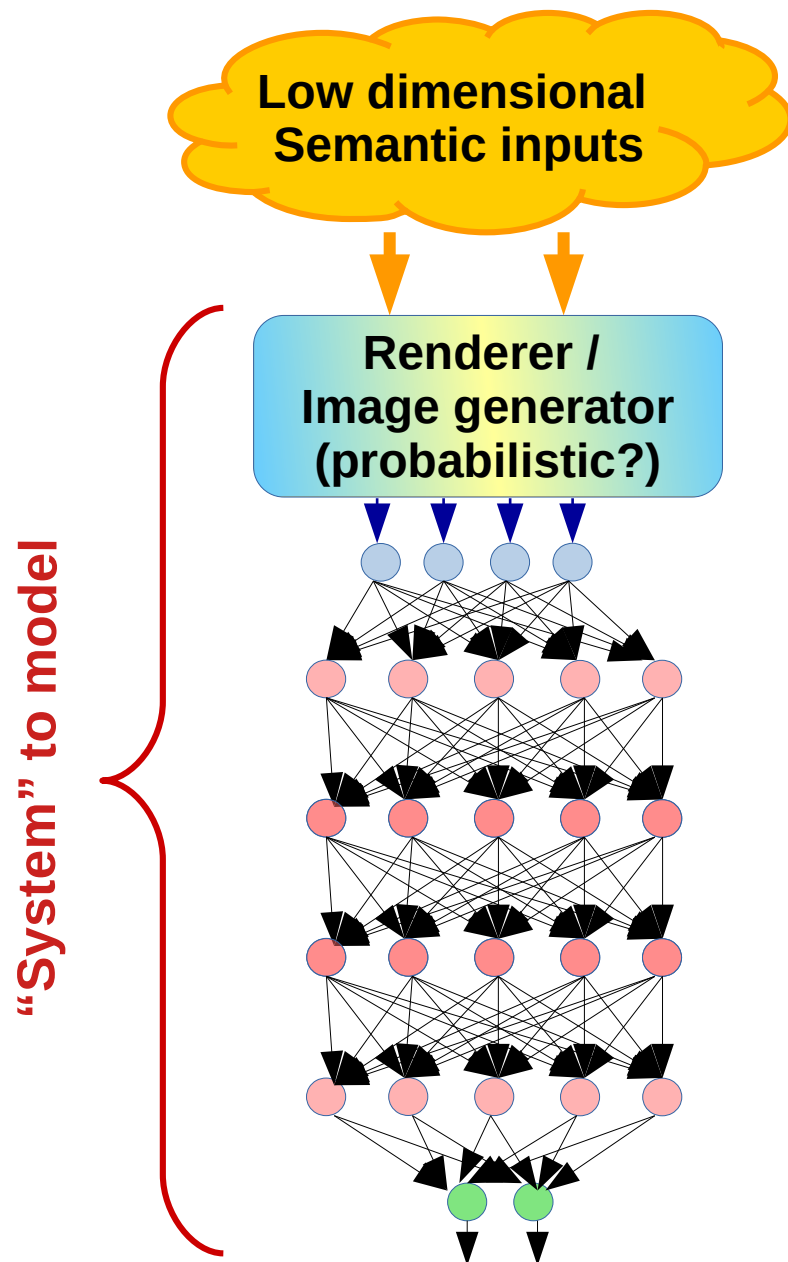
**Most images inconsequential, have no semantic similarity to what can possibly arise on a road**



“Too congested to accelerate”

**Can we restrict specs to a lower dimensional, smaller, meaningful input space?**

# Any Hope for Perceptual DNNs?



Time of Day: {Morning, Noon, Afternoon, Dusk, Night}  
Weather: {Clear, Cloudy, Snowing, Raining}  
Lanes: {Wide, Medium, Narrow, None}  
Road direction: {Straight, Bending}  
Other vehicles within 10m: {0, 1-3, 4-8, 9-15, > 15}  
Behaviour of other vehicles: {Lane disciplined, Chaotic}

Dimensions of semantic inp space = 6  
|Semantic inp space| =  $5 \times 4 \times 4 \times 2 \times 5 \times 2 = 1600$

Dimensions of image inp space =  $100 \times 100 \times 3 = 30000$   
|Image inp space| =  $256^{100 \times 100 \times 3}$

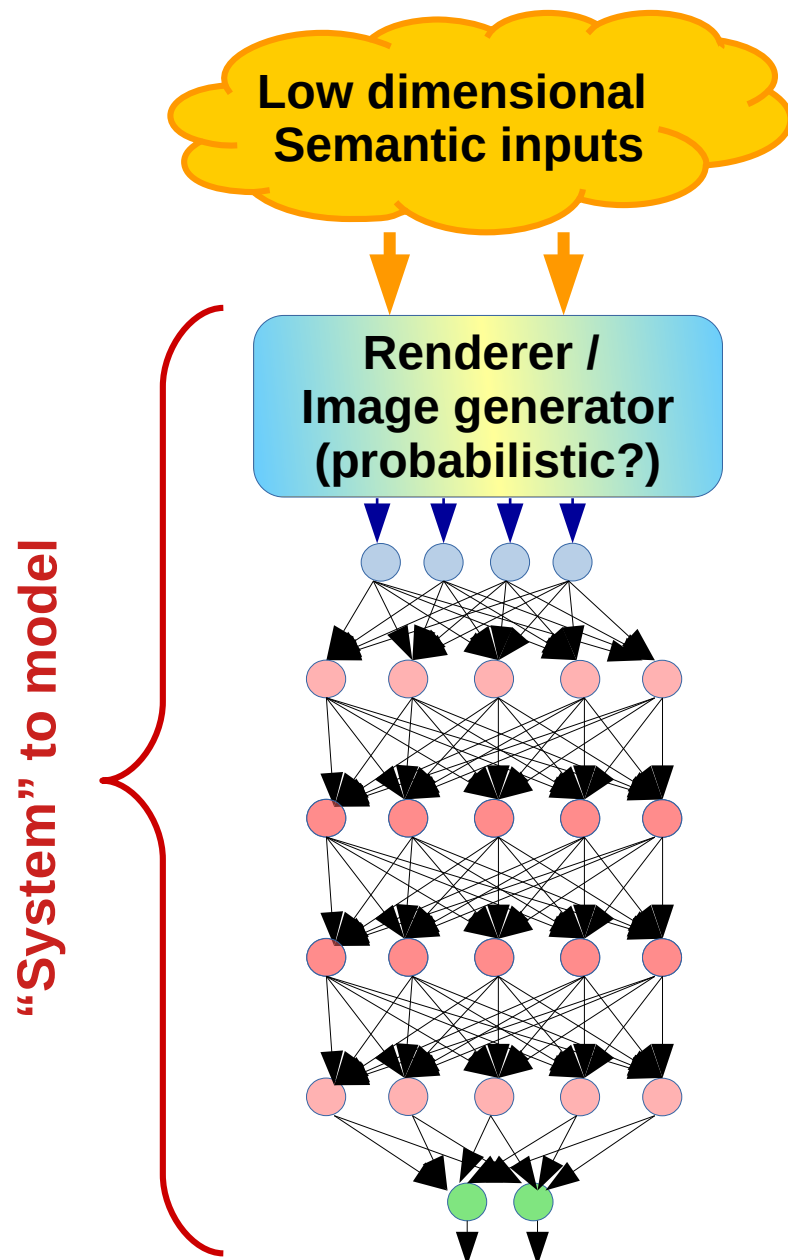
{ Pre-condition on semantic inputs  $\mathbf{s}$  }

$i \leftarrow \rho(\mathbf{s}); // \rho$ : Model of renderer

$y \leftarrow \nu(i); // \nu$ : Model of perceptual DNN

{ Post-condition on  $y$  }

# Any Hope for Perceptual DNNs?



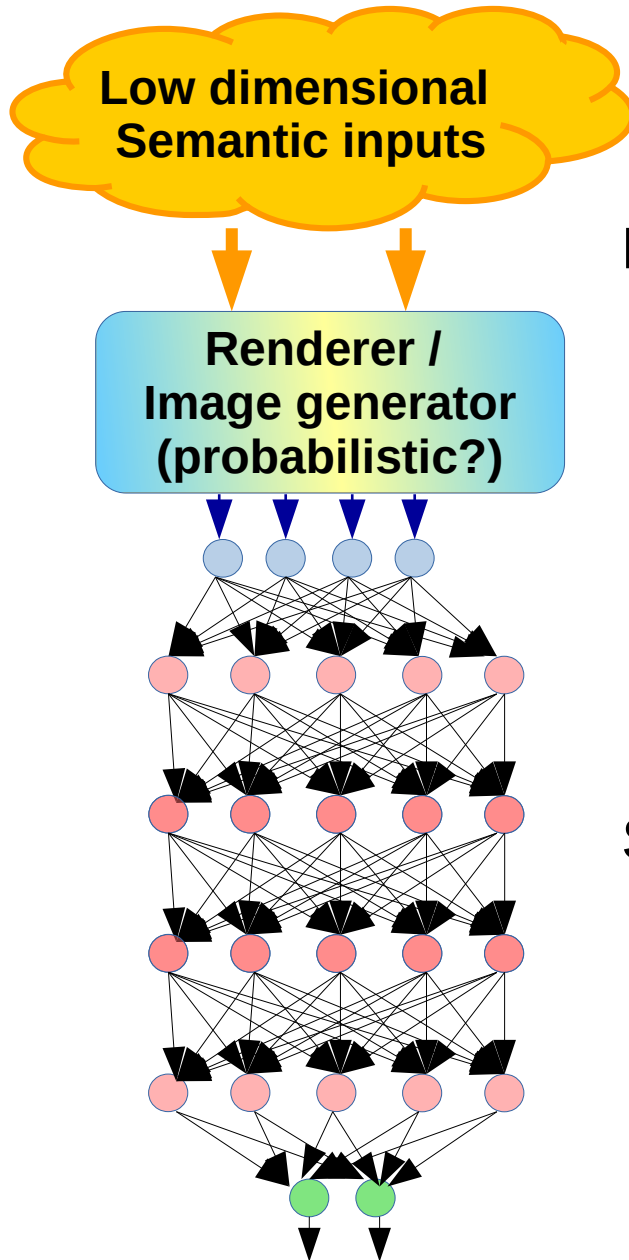
T: {Morning, Noon, Afternoon, Dusk, Night}  
W: {Clear, Cloudy, Snowing, Raining}  
L: {Wide, Medium, Narrow, None}  
Rd: {Straight, Bending}  
O: {0, 1-3, 4-8, 9-15, > 15}  
B: {Lane disciplined, Chaotic}

$$\{(O > 15) \vee (L = W) \wedge ((O \geq 9) \wedge (B = \text{Ch})) \vee (L = M) \wedge ((O \geq 9) \vee ((O \geq 4) \wedge (B = \text{Ch}))) \vee ((L = N) \vee (L = \text{None})) \wedge ((O \geq 4) \vee ((O \geq 1) \wedge (B = \text{Ch})))\}$$

$i \leftarrow \rho(T, W, L, Rd, O, B); // \rho$ : Model of renderer  
 $y \leftarrow \nu(i); // \nu$ : Model of perceptual DNN

{  $y$  = “Too congested to accelerate” }

# Any Hope for Perceptual DNNs?



Potential “**problems**”:

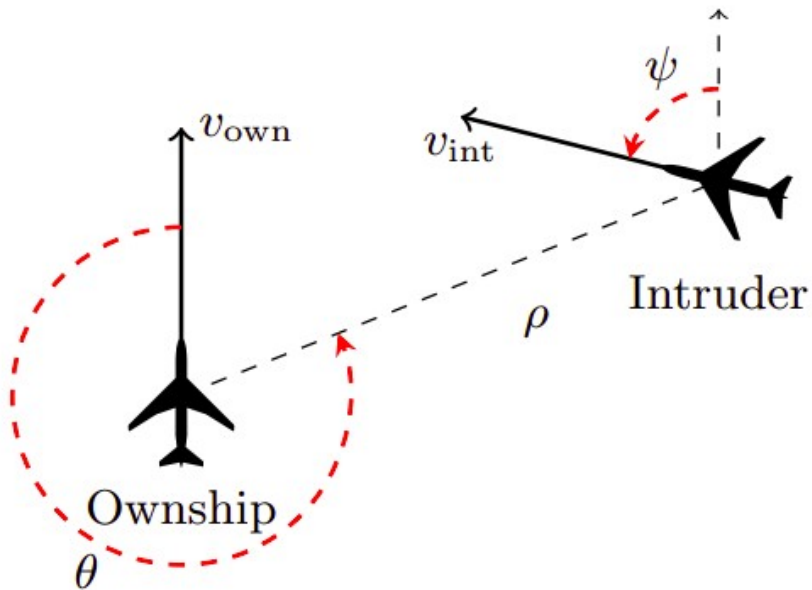
- Doesn’t cover entire input space
- Enrich semantic space to cover most/all meaningful inputs
- Use richer rendering modules
- **Need to model renderer**
- Use abstract / non-deterministic / probabilistic models

Significant “**benefits**”:

- Can eliminate large parts of irrelevant/meaningless input space
- Provide guarantees over large parts of meaningful input space

# One Spec vs Multiple Sub-specs

## ACAS-Xu



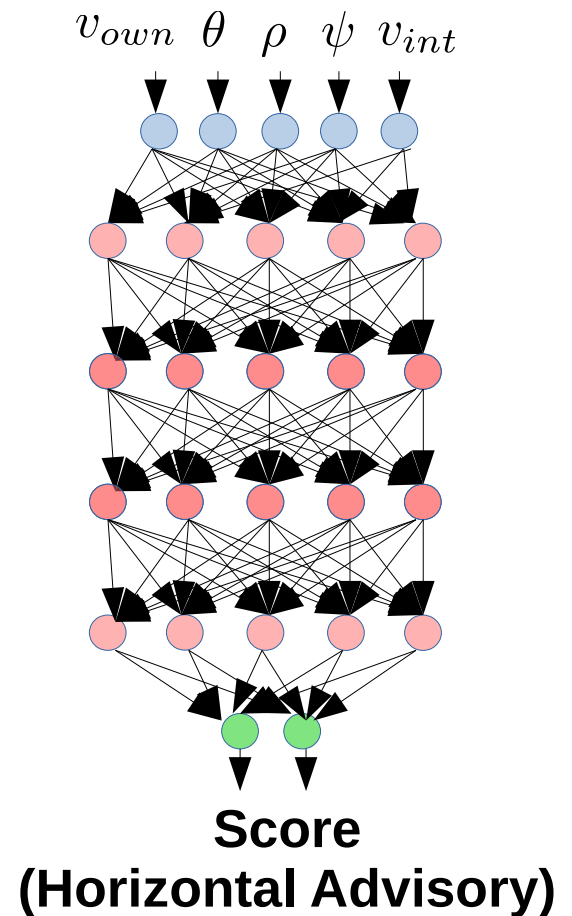
$$\{(\rho \geq 55947.691ft) \wedge (v_{own} \geq 1145ft/s) \wedge (v_{int} \leq 60ft/s)\}$$

$$\text{Score} \leftarrow \nu(\rho, v_{own}, v_{int}, \theta, \psi)$$

$$\{\text{Score}[\text{COC}] \leq 1500\}$$

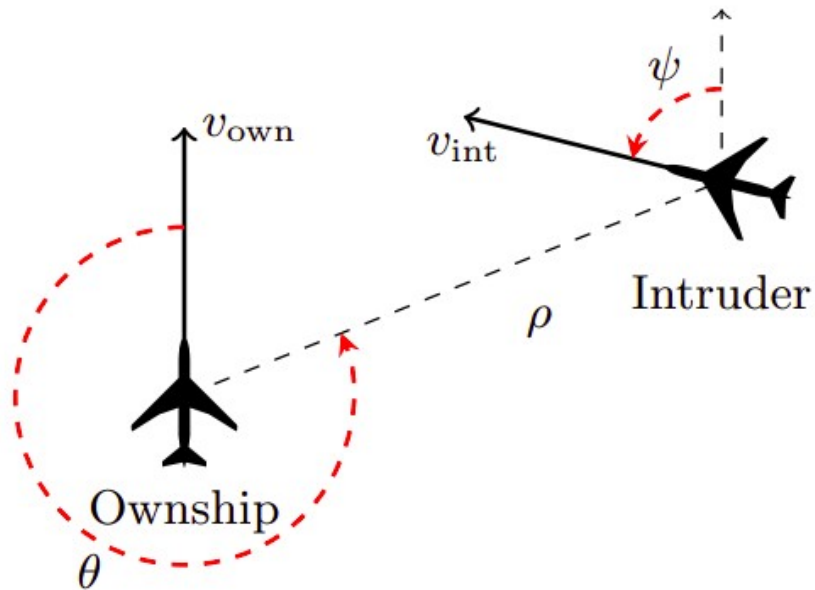
**Spec 1**

## Flight parameters



# One Spec vs Multiple Sub-specs

## ACAS-Xu



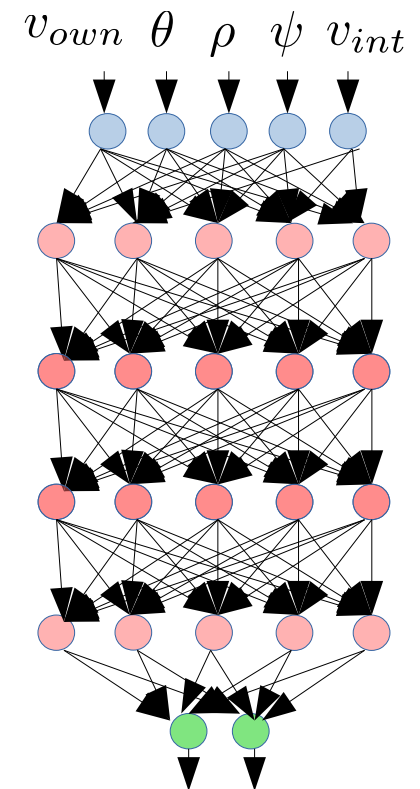
$$\{(0 \leq \rho \leq 60760ft) \wedge (1000 \leq v_{own} \leq 1200ft/s) \wedge (0 \leq v_{int} \leq 1200ft/s) \wedge (-3.141592 \leq \theta, \psi \leq 3.141592)\}$$

$$\mathbf{Score} \leftarrow \nu(\rho, v_{own}, v_{int}, \theta, \psi)$$

**Spec 7**

$$\{\operatorname{argmin}_x \mathbf{Score}[x] \notin \{\text{StrongRight}, \text{StrongLeft}\}\}$$

## Flight parameters



**Score**  
**(Horizontal Advisory)**

# One Spec vs Multiple Sub-specs

## ACAS-Xu

$$\{(v_{own} \geq 1000ft/s) \wedge (0 \leq v_{int} \leq 1200ft/s)\}$$

$$\text{Score} \leftarrow \nu(\rho, v_{own}, v_{int}, \theta, \psi)$$

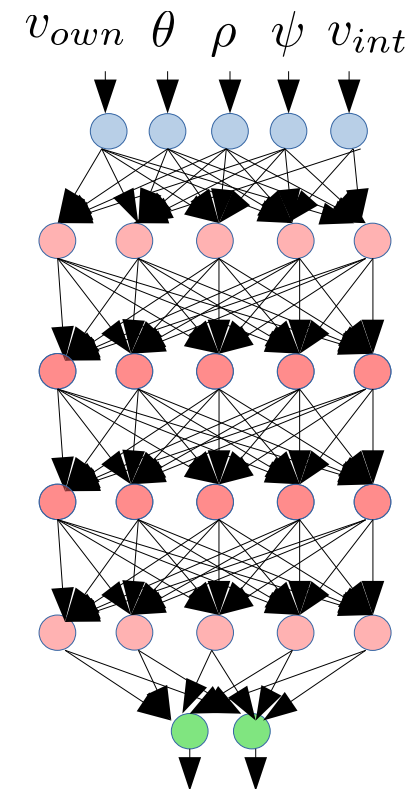
$$\left\{ \begin{array}{l} (\rho \geq 55947.691ft) \wedge (v_{own} \geq 1145ft/s) \wedge (v_{int} \leq 60ft/s) \\ \Rightarrow \text{Score}[\text{COC}] \leq 1500 \end{array} \right.$$

$\wedge$

$$\left\{ \begin{array}{l} (0 \leq \rho \leq 60760ft) \wedge (1000 \leq v_{own} \leq 1200ft/s) \wedge \\ (0 \leq v_{int} \leq 1200ft/s) \wedge (-3.141592 \leq \theta, \psi \leq 3.141592) \\ \Rightarrow \text{argmin}_x \text{Score}[x] \notin \{\text{StrongRight}, \text{StrongLeft}\} \end{array} \right.$$

**Specs 1+7**

Flight parameters



**Score**  
**(Horizontal Advisory)**



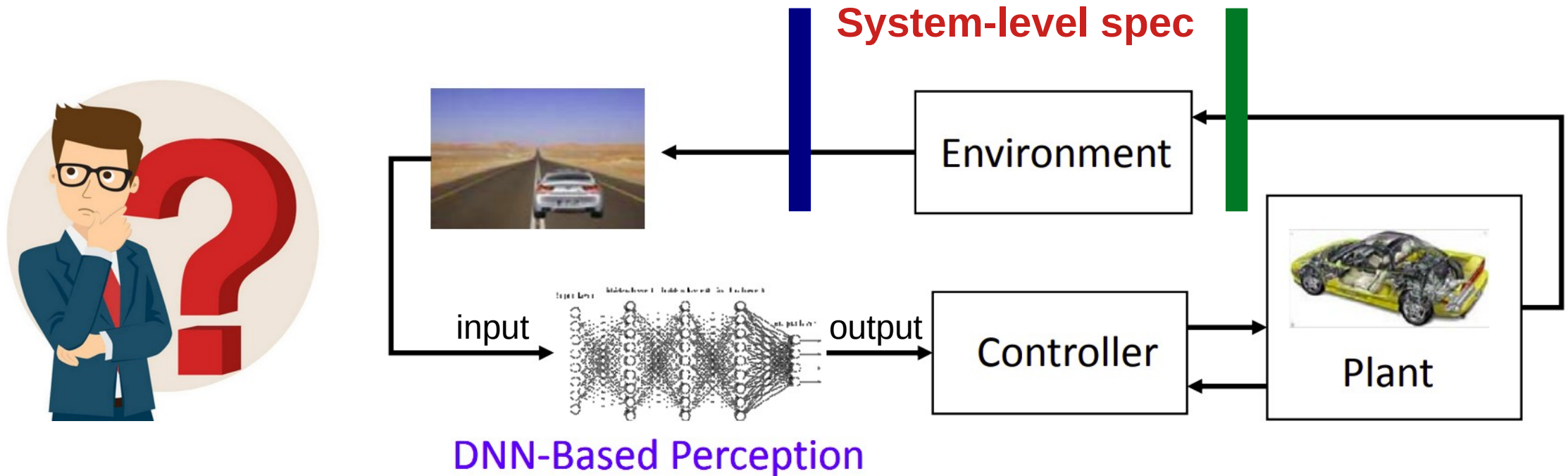
# One Spec vs Multiple Sub-specs

Multiple sub-specs generally preferred over one all-encompassing spec

- Separation of concerns
- Easy understandability
- Proofs often easier
- Modularly build spec over time



# Other Ways of Specifying Properties



Source: Seshia et al, Formal Verification of Deep Neural Networks, 2018

{ (own\_velocity > 30 km/h) and (road\_straight\_ahead) and (vehicles\_within\_5m = 0) }

Model of DNN + Controller + Plant

{ Steering = straight }

# Other Ways of Specifying Properties

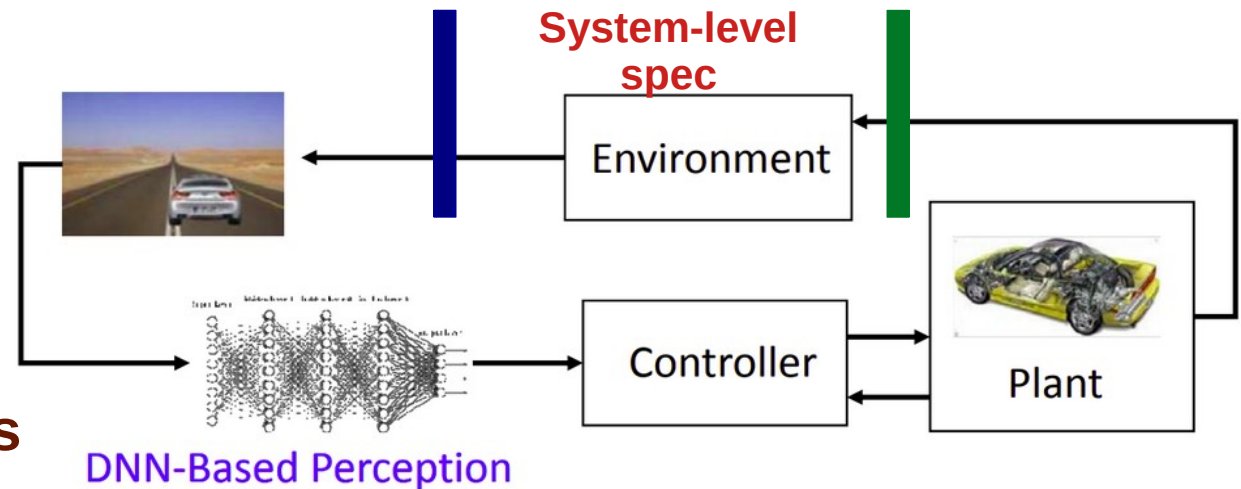
**No need for perceptual specs**

- Often easier to specify

**Require models of other components**

- May be harder to verify

**Classification errors of DNN may not translate to system level spec violations**



Source: Seshia et al, Formal Verification of Deep Neural Networks, 2018

{ (own\_velocity > 30 km/h) and (road\_straight\_ahead) and (vehicles\_within\_5m = 0) }

Model of DNN + Controller + Plant

{ Steering = straight }

# Specifying Properties of Neural Networks



**Pause n Reflect**

**DNNs are intended to mimic human reasoning**

**Is ideal human reasoning amenable to formal specification?**

**There are “boundaries” of acceptable/unacceptable human behaviour**

**Can we specify these boundaries?**

**Rules, laws, code of conduct**

**Do they have unique interpretations?**

**Do they evolve?**

**Is there a counterpart for neural networks?**