

# Are Uncertainty Capabilities of Evidential Deep Learning a Mirage?

Maohao Shen<sup>1\*</sup>, [J. Jon Ryu](#)<sup>1\*</sup>, Soumya Ghosh<sup>2</sup>,  
Yuheng Bu<sup>3</sup>, Prasanna Sattigeri<sup>2</sup>, Subhro Das<sup>2</sup>, Gregory W. Wornell<sup>1</sup>

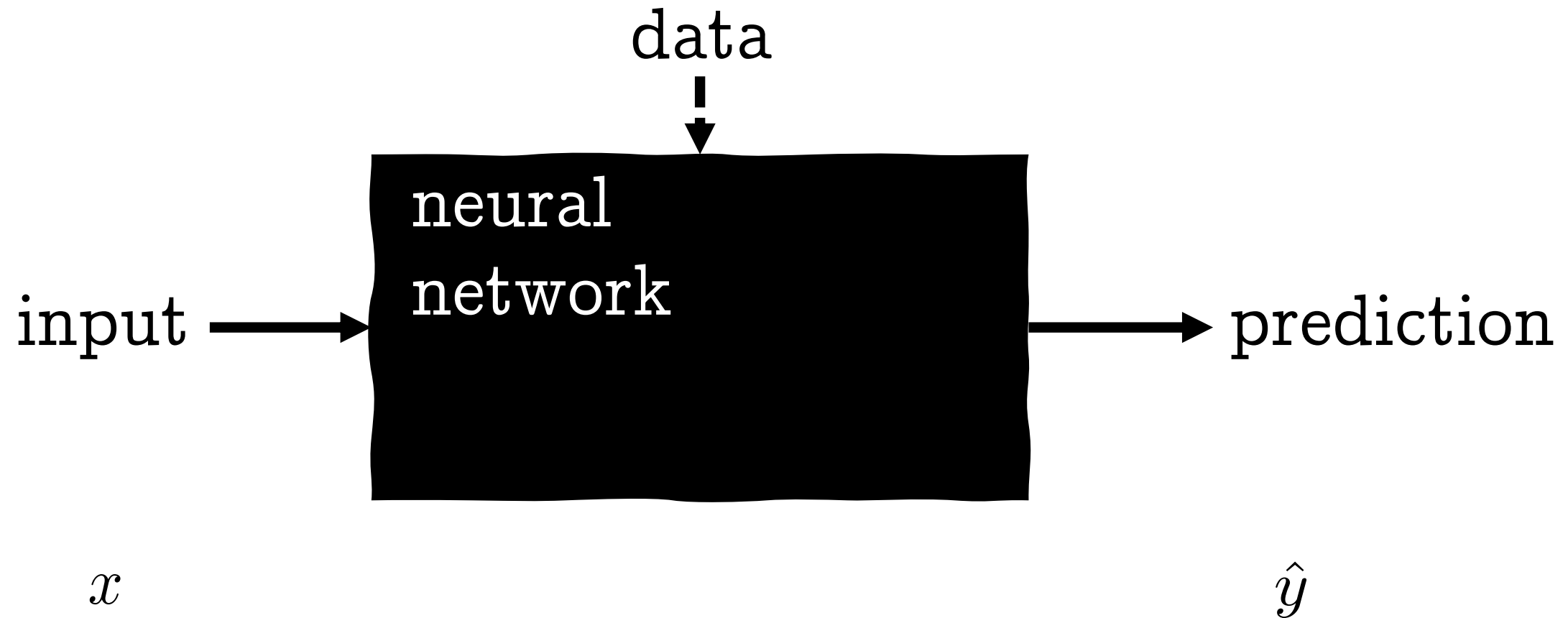
<sup>1</sup>MIT

<sup>2</sup>MIT-IBM Watson AI Lab, IBM Research

<sup>3</sup>University of Florida

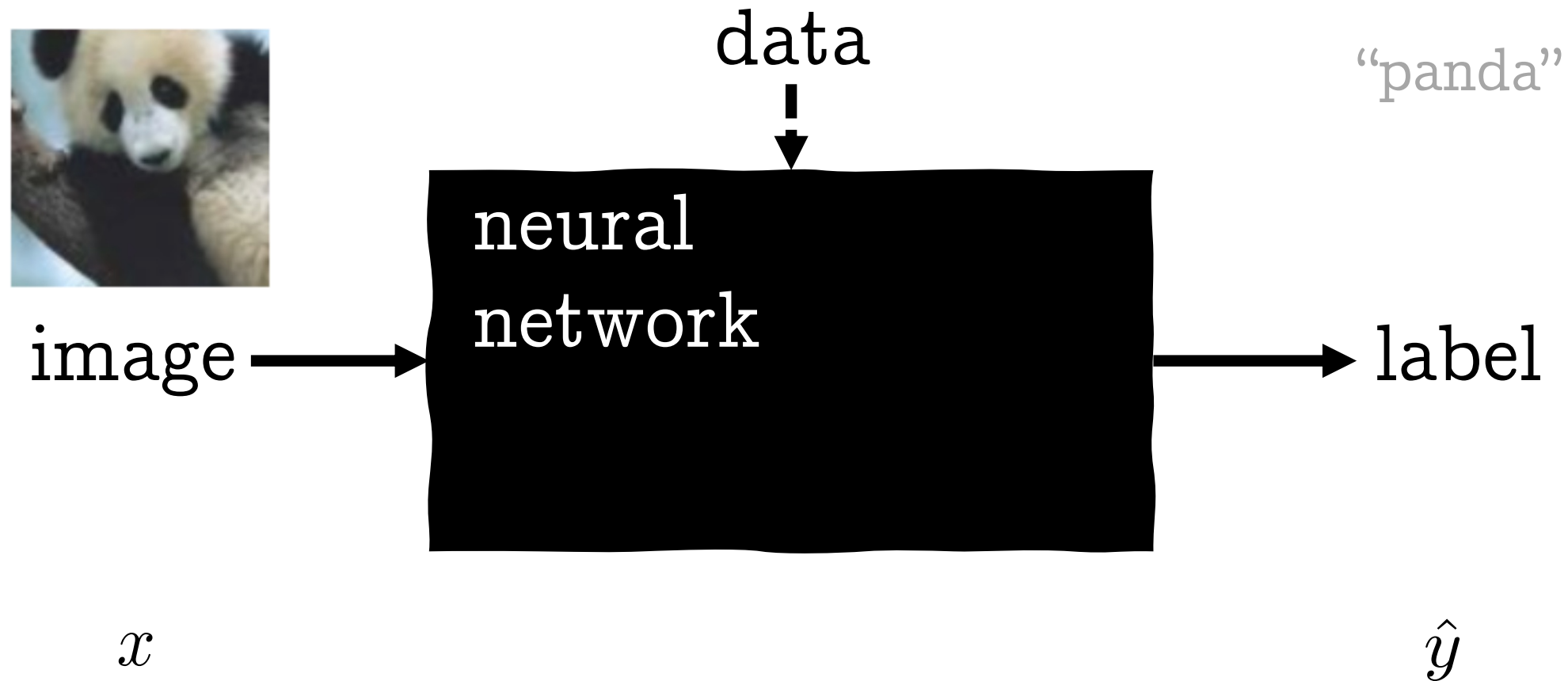


# Black-Box Prediction is Unreliable



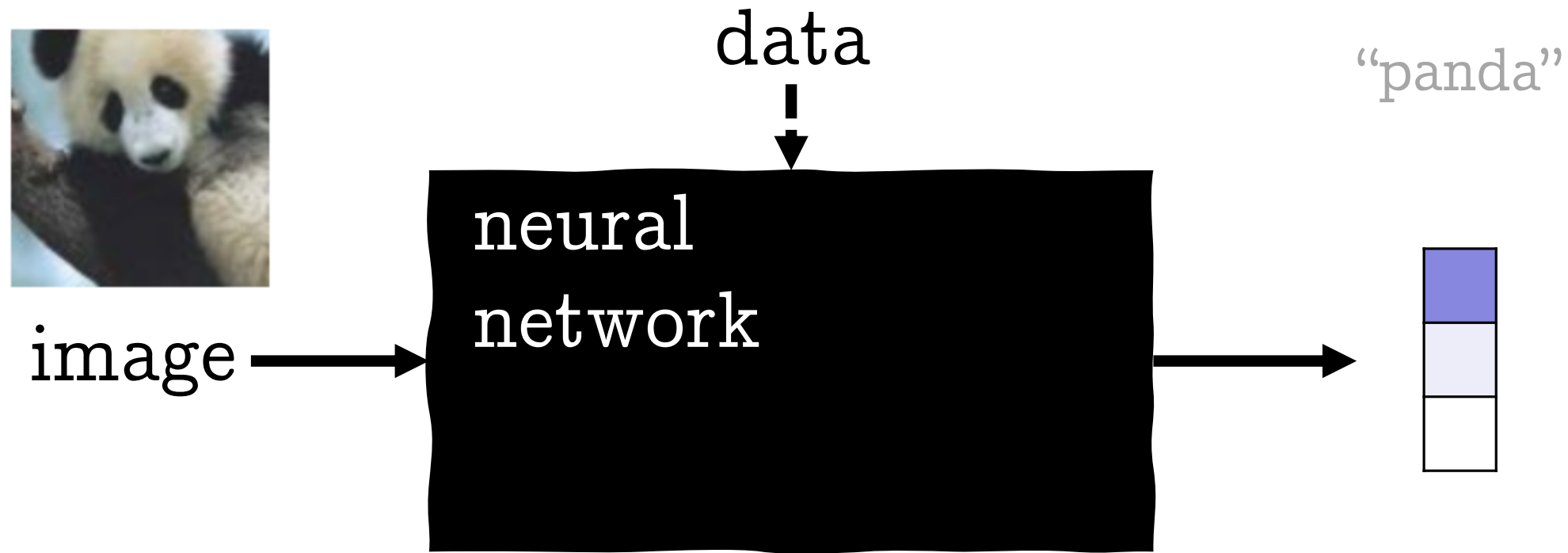
- Neural-network predictors are highly accurate...

# Black-Box Prediction is Unreliable



- Neural-network predictors are highly accurate...

# Black-Box Prediction is Unreliable

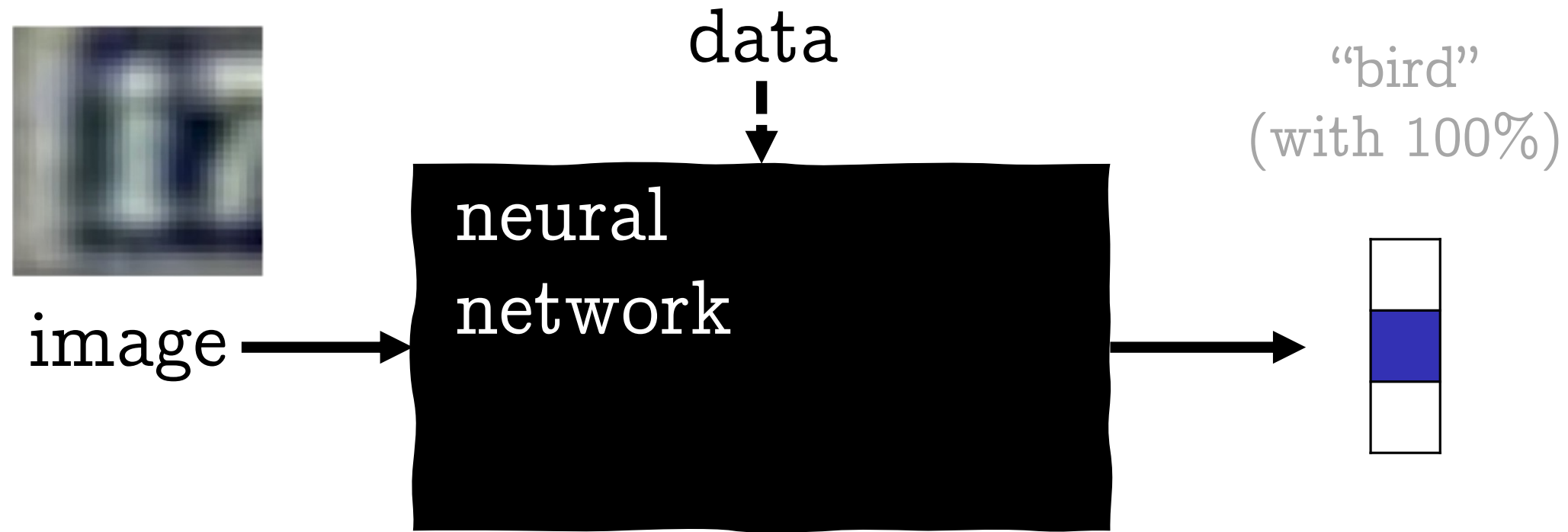


$x$

$$\pi_{\psi}(x) = (p_{\psi}(y|x))_{y=1}^C$$

- Neural-network predictors are highly accurate...

# Black-Box Prediction is Unreliable

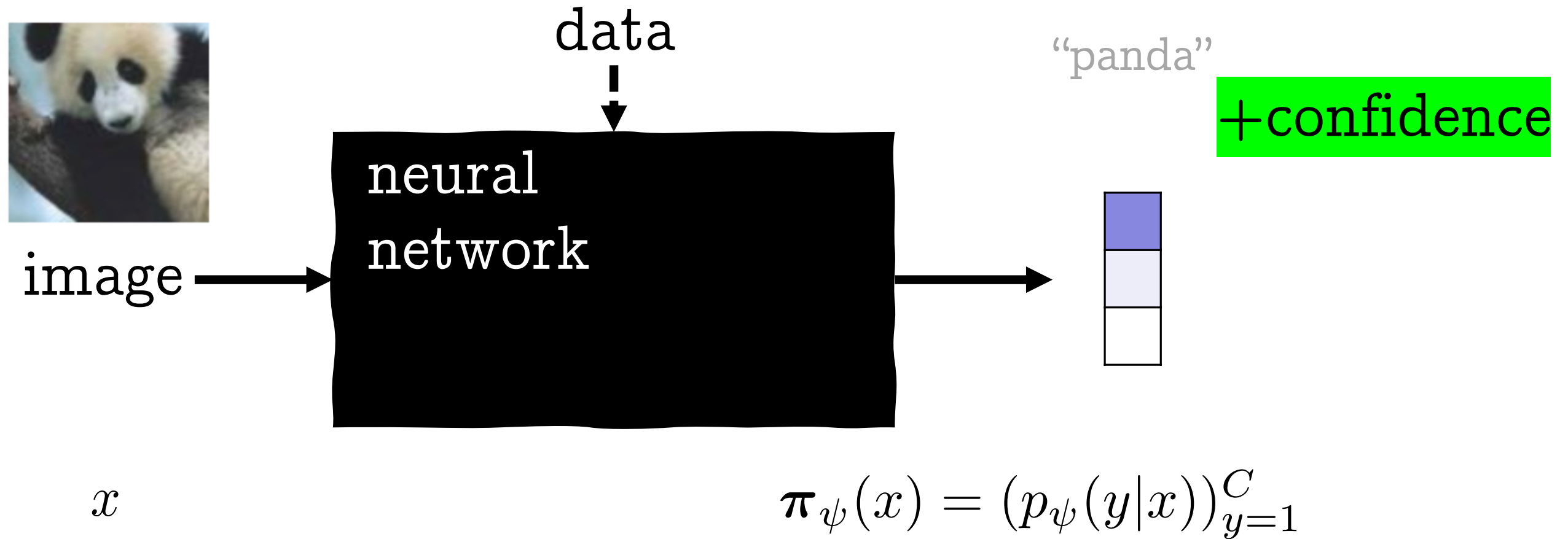


$x$

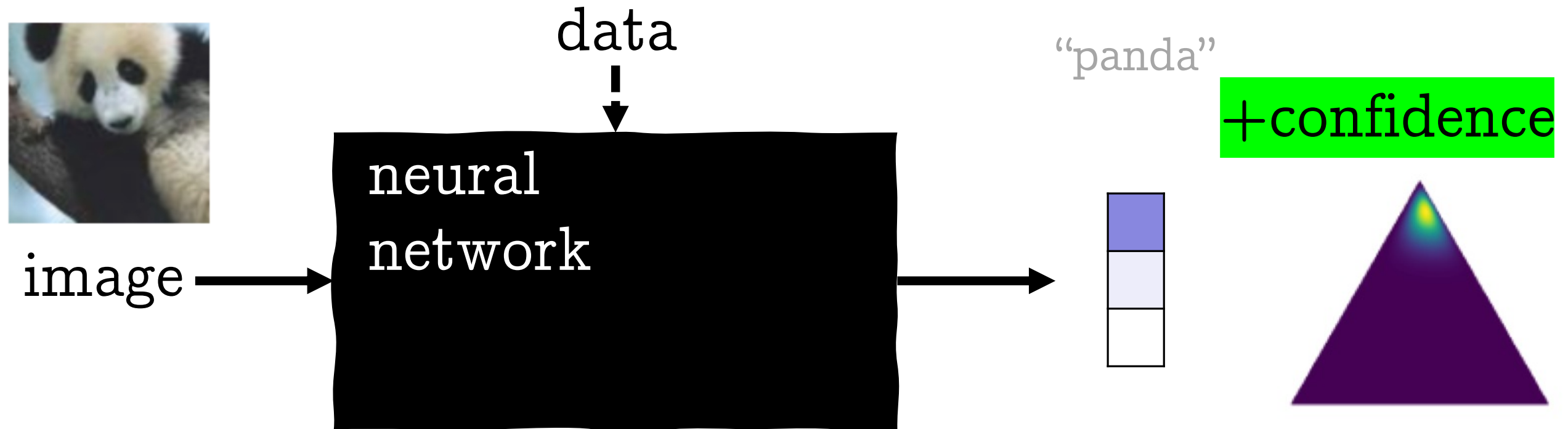
$$\pi_{\psi}(x) = (p_{\psi}(y|x))_{y=1}^C$$

- Neural-network predictors are highly accurate...
- But often **overconfident** for out-of-distribution (OOD) data

# We Need Uncertainty Quantification!



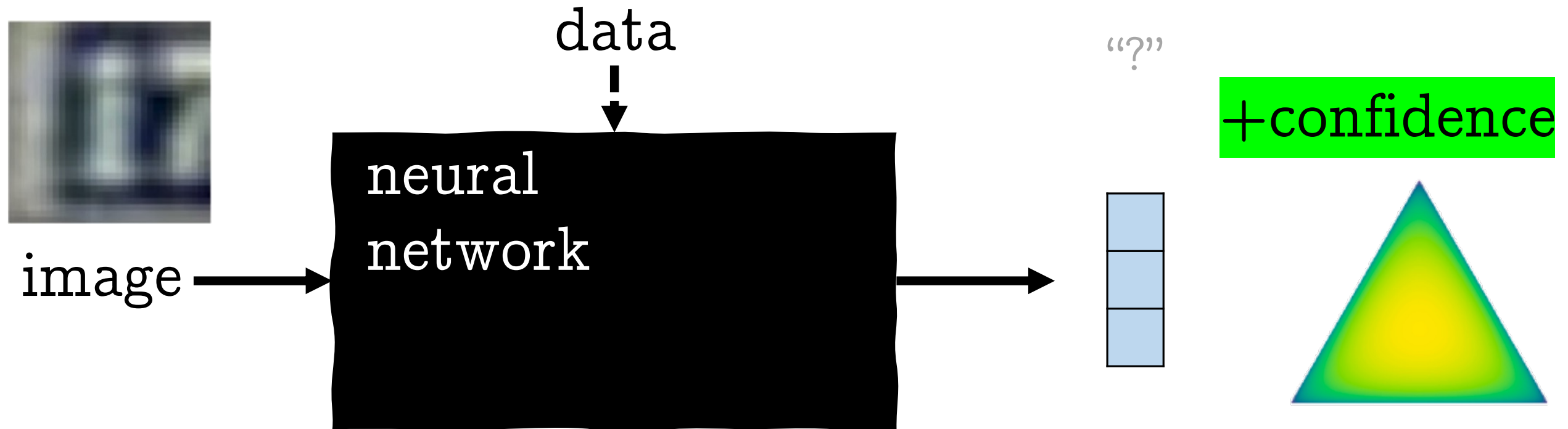
# We Need Uncertainty Quantification!



$x$

$$\pi_{\psi}(x) = (p_{\psi}(y|x))_{y=1}^C \quad p_{\psi}(\pi|x)$$

# We Need Uncertainty Quantification!

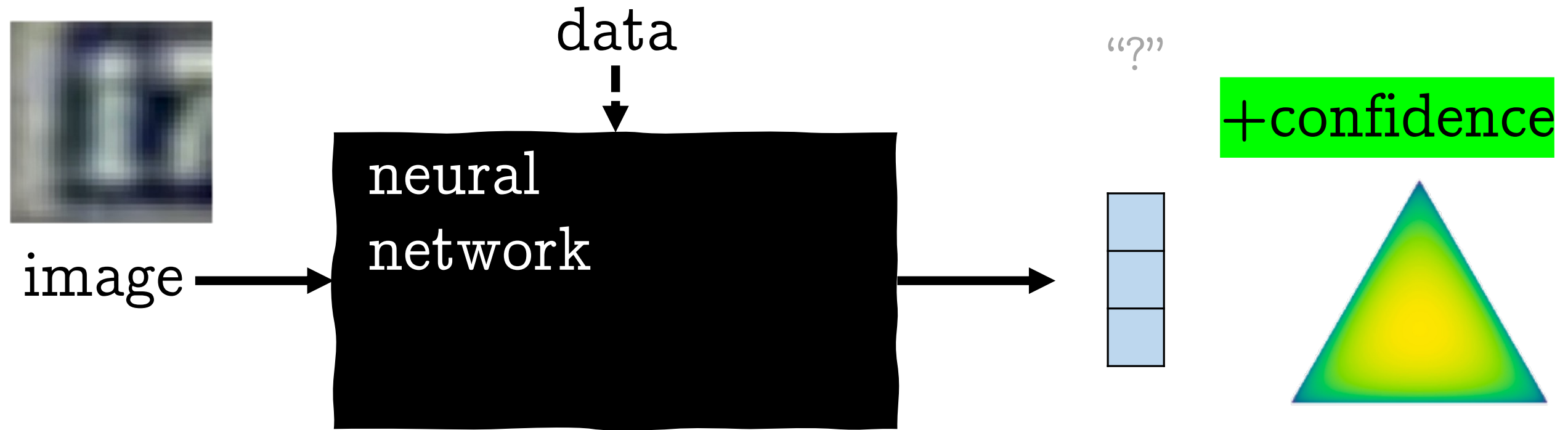


$x$

$$\boldsymbol{\pi}_{\psi}(x) = (p_{\psi}(y|x))_{y=1}^C \quad p_{\psi}(\boldsymbol{\pi}|x)$$



# We Need Uncertainty Quantification!



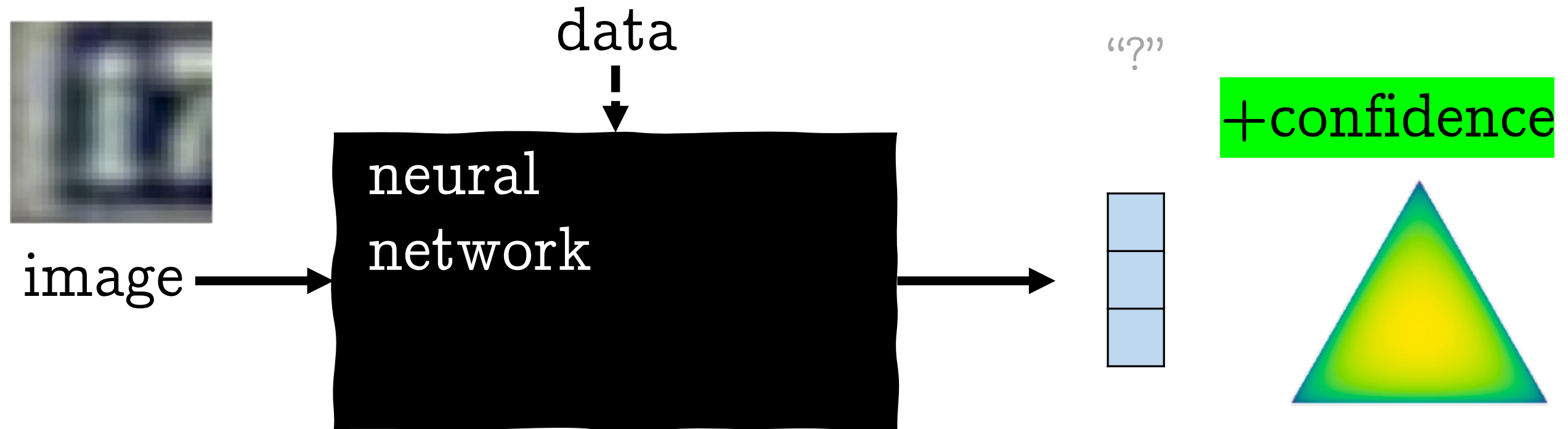
$x$

$$\pi_{\psi}(x) = (p_{\psi}(y|x))_{y=1}^C \quad p_{\psi}(\pi|x)$$

Different ways of inducing  $p_{\psi}(\pi|x)$

- Bayesian methods: variational inference, MCMC, Monte Carlo Dropout, ...
- Frequentist methods: jackknife, bootstrap, ...
- Ensemble methods

# We Need Uncertainty Quantification!



$x$

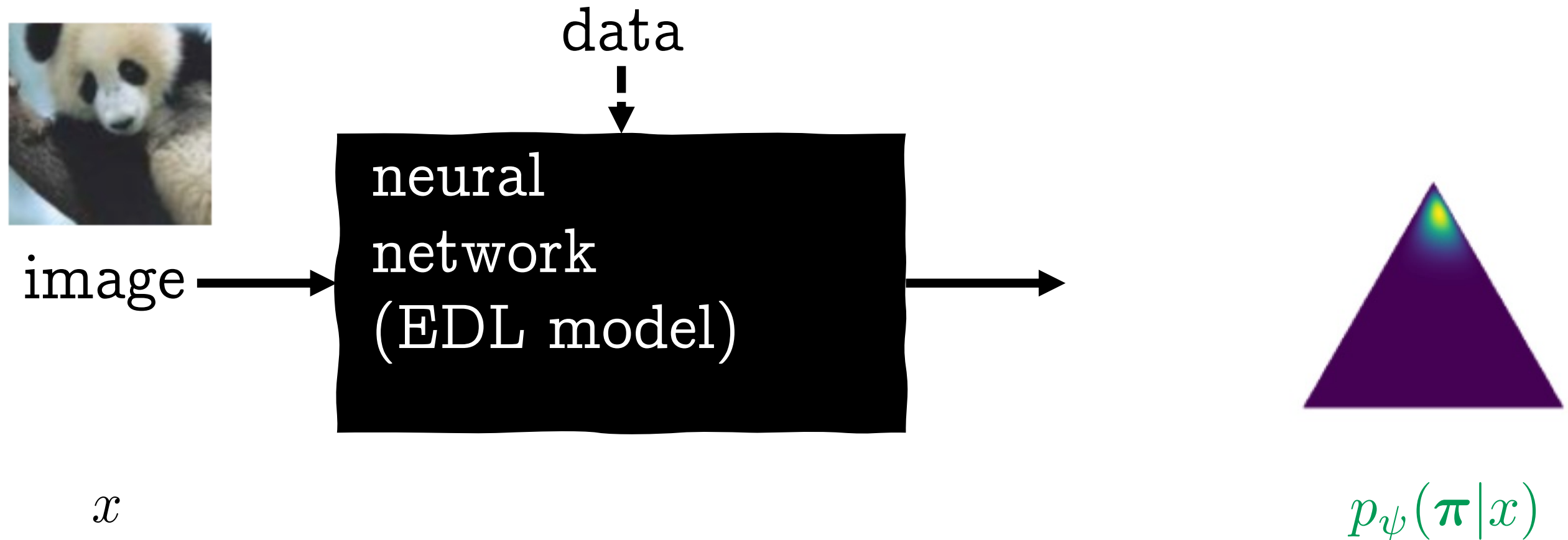
$$\pi_{\psi}(x) = (p_{\psi}(y|x))_{y=1}^C \quad p_{\psi}(\pi|x)$$

Different ways of inducing  $p_{\psi}(\pi|x)$

- Bayesian methods: variational inference, MCMC, Monte Carlo Dropout, ...
- Frequentist methods: jackknife, bootstrap, ...
- Ensemble methods

**Computationally Inefficient!**

# Evidential Deep Learning (EDL)



- Directly train a **single neural network** that outputs  $p_{\psi}(\pi|x)$
- Empirical successes for downstream tasks (e.g., OOD detection)
- Lack of theoretical understanding
- Recent works have reported spurious behaviors

# Demystifying EDL Methods

- Q1. What do EDL methods learn as **uncertainty**?
- Q2. **Why** are the EDL methods **successful**?
- Q3. How can we **make** EDL methods **more reliable**?

# Unifying EDL Objectives: A New Taxonomy

---

Method (name of loss)
-----------------------

---

FPriorNet (F-KL loss) <a href="#">[7]</a>
---

RPriorNet (R-KL loss) <a href="#">[8]</a>
---

EDL (MSE loss) <a href="#">[10]</a>
-------------------------------------

Belief Matching (VI loss) <a href="#">[12, 13]</a>
--

PostNet (UCE loss) <a href="#">[15]</a>
---

NatPN (UCE loss) <a href="#">[20]</a>
---------------------------------------

---

What's the common principle  
behind all these objectives?

# Unifying EDL Objectives: A New Taxonomy

---

Method (name of loss)

---

FPriorNet (F-KL loss) [7]

RPriorNet (R-KL loss) [8]

EDL (MSE loss) [10]

Belief Matching (VI loss) [12, 13]

PostNet (UCE loss) [15]

NatPN (UCE loss) [20]

---

What's the common principle  
behind all these objectives?

$$\mathcal{L}(\psi) := \mathbb{E}_{p(x,y)} [D(\underbrace{p^{(\nu)}(\boldsymbol{\pi} | y)}_{\text{"fixed" uncertainty target}}, \underbrace{p_{\psi}(\boldsymbol{\pi} | x)}_{\text{EDL model}})] + \gamma_{\text{ood}} \mathbb{E}_{p_{\text{ood}}(x)} [D(p(\boldsymbol{\pi}), \underbrace{p_{\psi}(\boldsymbol{\pi} | x)}_{\text{EDL model}})]$$

$\underbrace{\hspace{15em}}_{\text{in-distribution objective}} \qquad \underbrace{\hspace{15em}}_{\text{OOD objective}}$

# Unifying EDL Objectives: A New Taxonomy

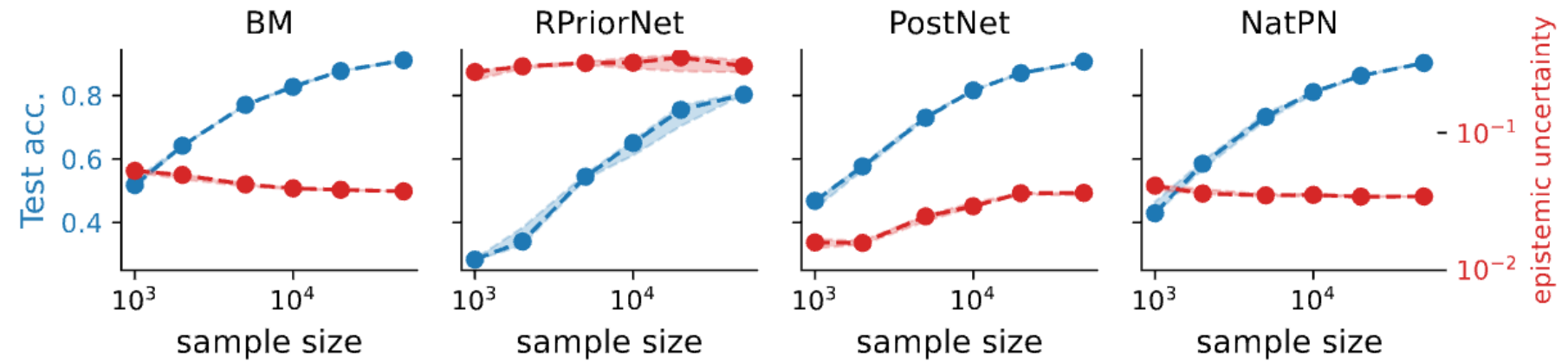
Method (name of loss)	likelihood	$D(\cdot, \cdot)$	prior $\alpha_0$	$\gamma_{\text{ood}}$	$\alpha_\psi(x)$ parameterization
FPriorNet (F-KL loss) [7]	categorical	fwd. KL	$= \mathbb{1}_C$	$> 0$	direct
RPriorNet (R-KL loss) [8]	categorical	rev. KL	$= \mathbb{1}_C$	$> 0$	direct
EDL (MSE loss) [10]	Gaussian	rev. KL	$= \mathbb{1}_C$	$= 0$	direct
Belief Matching (VI loss) [12, 13]	categorical	rev. KL	$\in \mathbb{R}_{>0}^C$	$= 0$	direct
PostNet (UCE loss) [15]	categorical	rev. KL	$= \mathbb{1}_C$	$= 0$	density w/ single flow
NatPN (UCE loss) [20]	categorical	rev. KL	$= \mathbb{1}_C$	$= 0$	density w/ multiple flows

$$\mathcal{L}(\psi) := \mathbb{E}_{p(x,y)} [D(\underbrace{p^{(\nu)}(\pi|y)}_{\text{"fixed" uncertainty target}}, \underbrace{p_\psi(\pi|x)}_{\text{EDL model}})] + \gamma_{\text{ood}} \mathbb{E}_{p_{\text{ood}}(x)} [D(p(\pi), \underbrace{p_\psi(\pi|x)}_{\text{EDL model}})]$$

in-distribution objective
OOD objective

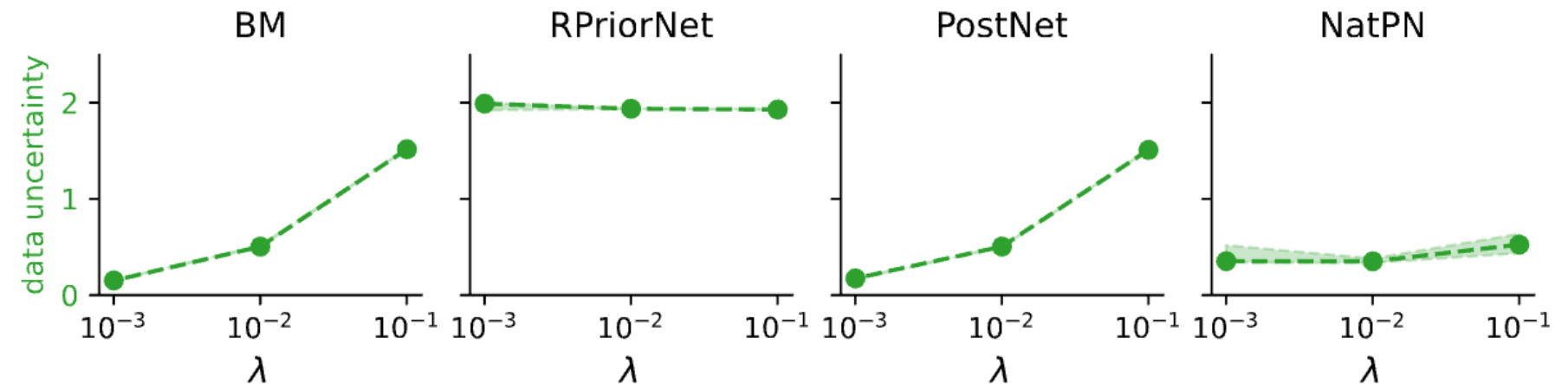
# Explain Spurious Behaviors of Learned Uncertainties

Non-vanishing  
learned epistemic uncertainty



(a) Epistemic Uncertainty: CIFAR10

Hyperparameter sensitive  
learned aleatoric uncertainty



(b) Aleatoric Uncertainty: CIFAR10



# Explain Success of EDL Methods

$$\mathcal{L}(\psi) := \mathbb{E}_{p(x,y)} [D(\underbrace{p^{(\nu)}(\pi|y)}_{\text{"fixed" uncertainty target}}, \underbrace{p_\psi(\pi|x)}_{\text{EDL model}})] + \gamma_{\text{ood}} \mathbb{E}_{p_{\text{ood}}(x)} [D(p(\pi), \underbrace{p_\psi(\pi|x)}_{\text{EDL model}})]$$

in-distribution objective
OOD objective

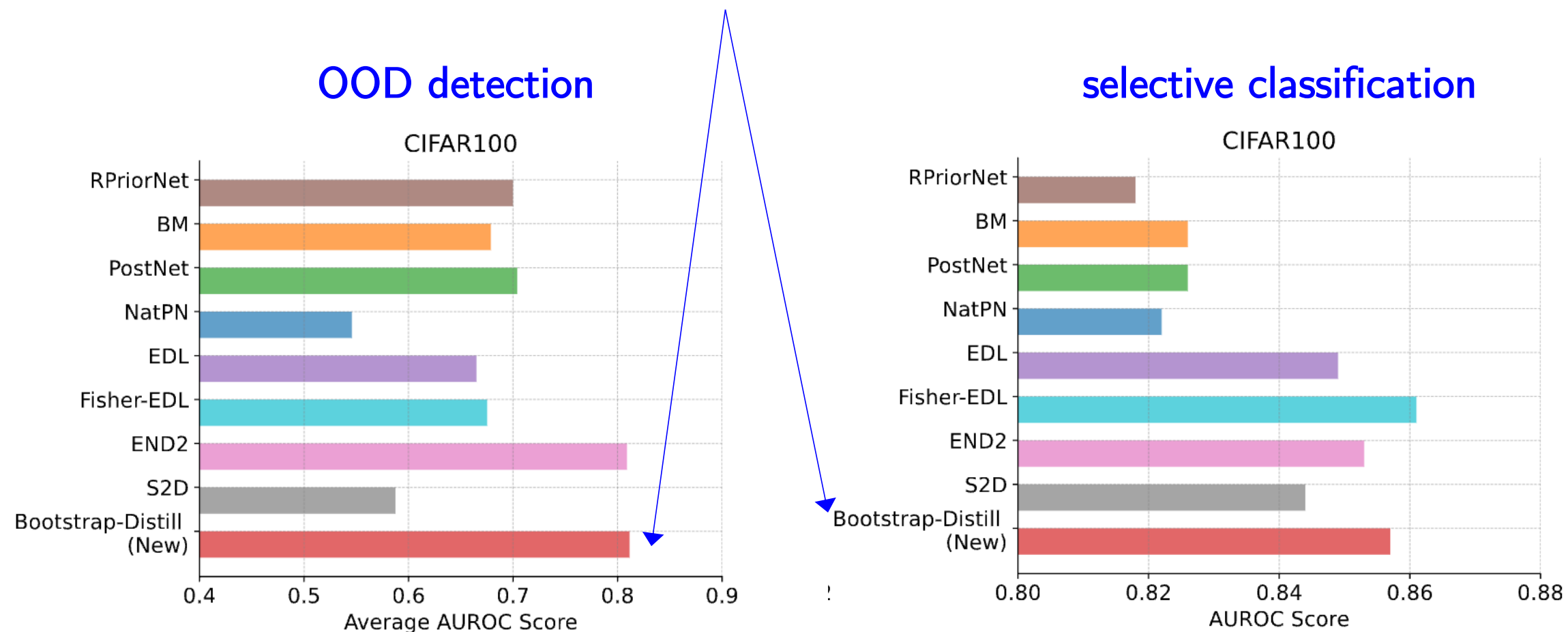
$\approx$

$$\underbrace{-\mathbb{E}_{p(x,y)} [\log p_\psi(y|x)] + \tau \{ \mathbb{E}_{p(x)} [\max(0, E_\phi(x) - m_{\text{id}})^2] }_{\text{in-distribution objective}} + \underbrace{\mathbb{E}_{p_{\text{ood}}(x)} [\max(0, m_{\text{ood}} - E_\phi(x))^2] }_{\text{OOD objective}} \},$$

- EDL methods can be better understood as **EBM-based OOD detector**

# How to Improve EDL?

- All the issues of EDL arise from **the ignorance of model stochasticity**
- **Conjecture:** **Incorporating external stochasticity** is the key for meaningful UQ and EDL should be used for “**distillation**” for fast inference
- Show **distilling** randomness in **Bootstrap** can achieve SOTA performance



---

# Are Uncertainty Quantification Capabilities of Evidential Deep Learning a Mirage?

---

**Maohao Shen<sup>1,\*</sup>, J. Jon Ryu<sup>1,\*</sup>, Soumya Ghosh<sup>2,†</sup>,  
Yuheng Bu<sup>3</sup>, Prasanna Sattigeri<sup>2</sup>, Subhro Das<sup>2</sup>, Gregory W. Wornell<sup>1</sup>**

<sup>1</sup>Department of EECS, MIT, Cambridge, MA 02139

<sup>2</sup>MIT-IBM Watson AI Lab, IBM Research, Cambridge, MA 02142

<sup>3</sup>Department of ECE, University of Florida, Gainesville, FL 32611

[{maohao,jongha,gww}@mit.edu](mailto:{maohao,jongha,gww}@mit.edu),  
[{ghoshoso,prasanna}@us.ibm.com](mailto:{ghoshoso,prasanna}@us.ibm.com), [subhro.das@ibm.com](mailto:subhro.das@ibm.com),  
[buyuheng@ufl.edu](mailto:buyuheng@ufl.edu)