

Deep k-Nearest Neighbors

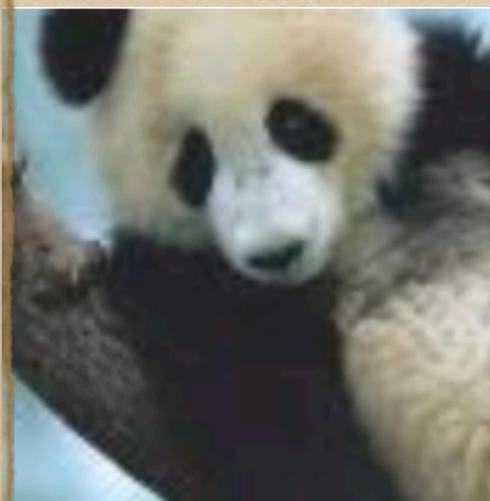
Buhua Liu

DkNN

- ◆ Introduction
- ◆ Problem Definition
- ◆ Methodology
- ◆ Experiments
- ◆ Conclusion

DkNN

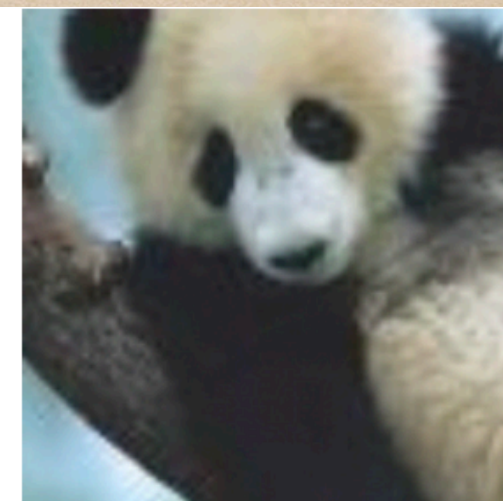
- ◆ Introduction
- ◆ Problem Definition
- ◆ Methodology
- ◆ Experiments
- ◆ Conclusion



+ .007 ×



=



x

“panda”

57.7% confidence

$\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”

8.2% confidence

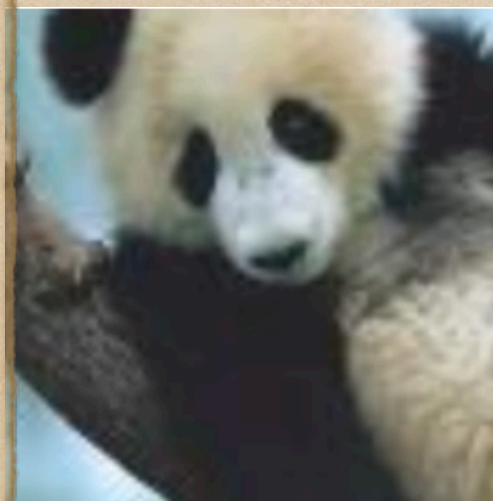
$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$

“gibbon”

99.3 % confidence

Adversarial Examples

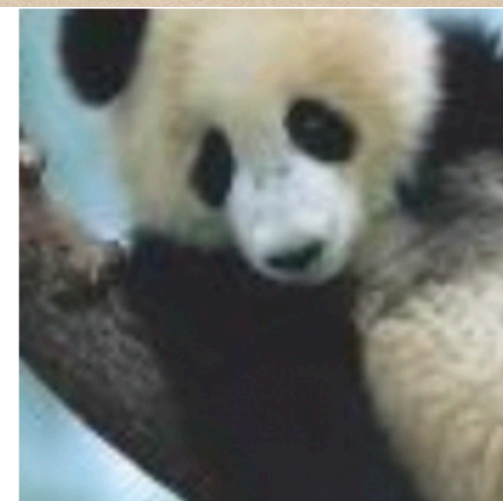
<http://arxiv.org/abs/1412.6572>

 x

“panda”
57.7% confidence

 $+ .007 \times$  $\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”
8.2% confidence

 $=$  $x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$

“gibbon”
99.3 % confidence

Adversarial Examples

<http://arxiv.org/abs/1412.6572>

DkNN

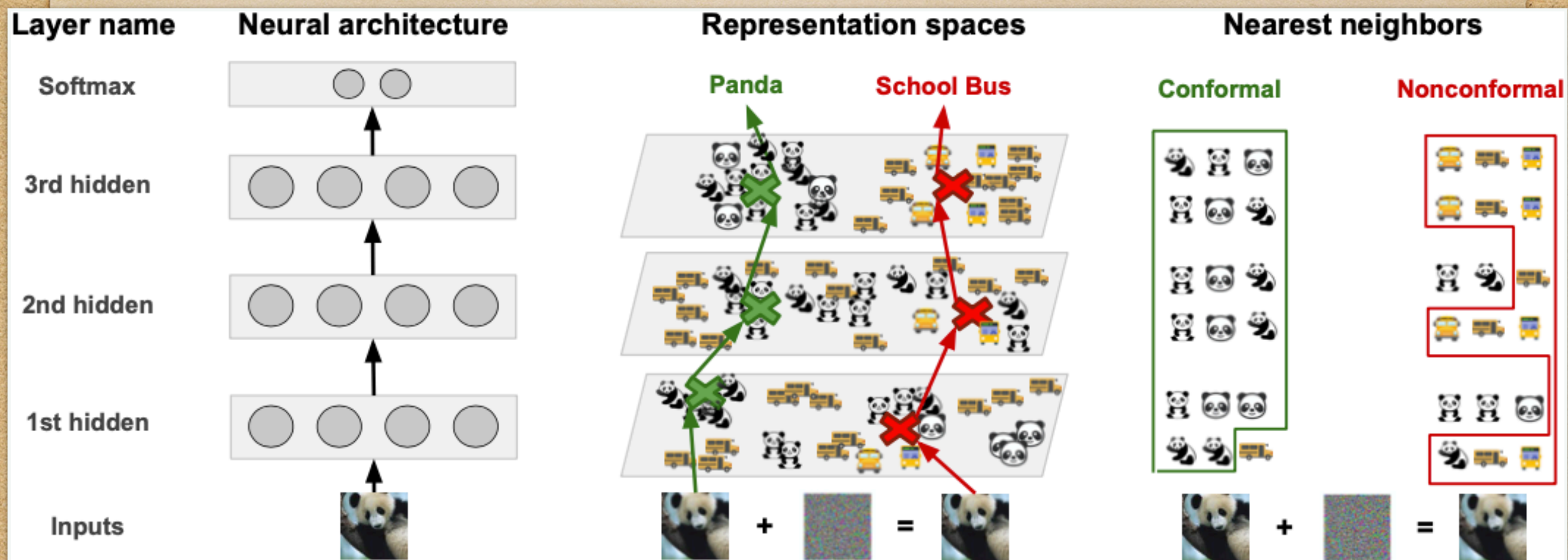
- ◆ Introduction
- ◆ Problem Definition
- ◆ Methodology
- ◆ Experiments
- ◆ Conclusion

What is wrong?

- ◆ Misclassification—Lack of robustness
- ◆ Unreliable confidence estimates
- ◆ Lack of interpretability—Over-parameterized DNN

DkNN

- ◆ Introduction
- ◆ Problem Definition
- ◆ Methodology
- ◆ Experiments
- ◆ Conclusion



Intuition behind the DkNN

<http://arxiv.org/abs/1803.04765>

Algorithm 1 – Deep k-Nearest Neighbor.

Input: training data (X, Y) , calibration data (X^c, Y^c)

Input: trained neural network f with l layers

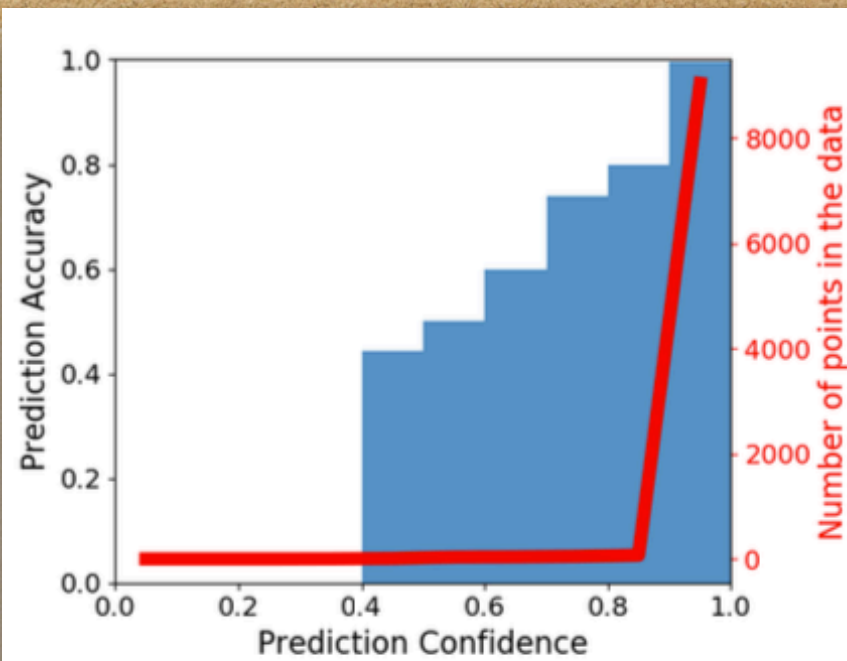
Input: number k of neighbors

Input: test input z

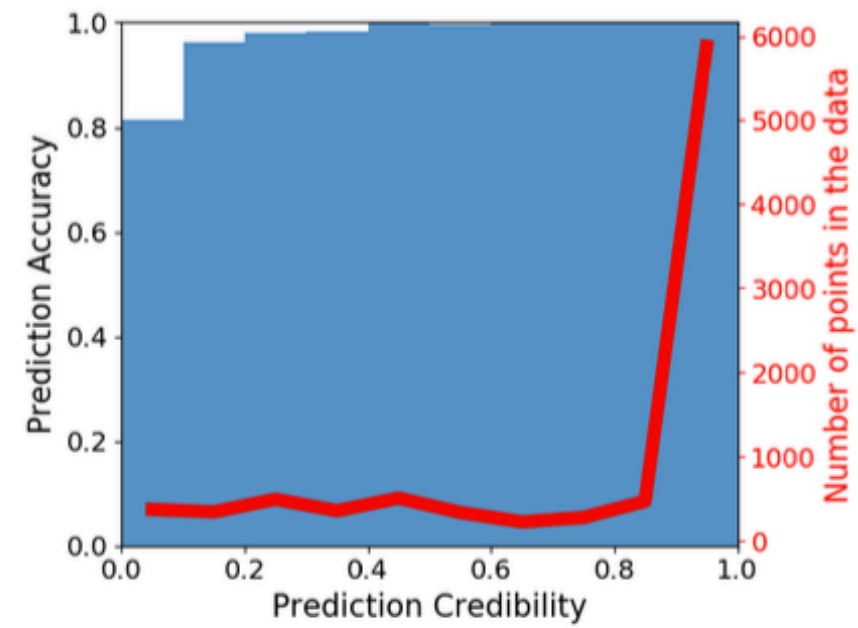
- 1: // Compute layer-wise k nearest neighbors for test input z
 - 2: **for** each layer $\lambda \in 1..l$ **do**
 - 3: $\Gamma \leftarrow k$ points in X closest to z found w/ LSH tables
 - 4: $\Omega_\lambda \leftarrow \{Y_i : i \in \Gamma\}$ \triangleright Labels of k inputs found
 - 5: **end for**
 - 6: // Compute prediction, confidence and credibility
 - 7: $A = \{\alpha(x, y) : (x, y) \in (X^c, Y^c)\}$ \triangleright Calibration
 - 8: **for** each label $j \in 1..n$ **do**
 - 9: $\alpha(z, j) \leftarrow \sum_{\lambda \in 1..l} |i \in \Omega_\lambda : i \neq j|$ \triangleright Nonconformity
 - 10: $p_j(z) = \frac{|\{\alpha \in A : \alpha \geq \alpha(z, j)\}|}{|A|}$ \triangleright empirical p -value
 - 11: **end for**
 - 12: prediction $\leftarrow \arg \max_{j \in 1..n} p_j(z)$
 - 13: confidence $\leftarrow 1 - \max_{j \in 1..n, j \neq \text{prediction}} p_j(z)$
 - 14: credibility $\leftarrow \max_{j \in 1..n} p_j(z)$
 - 15: **return** prediction, confidence, credibility
-

DkNN

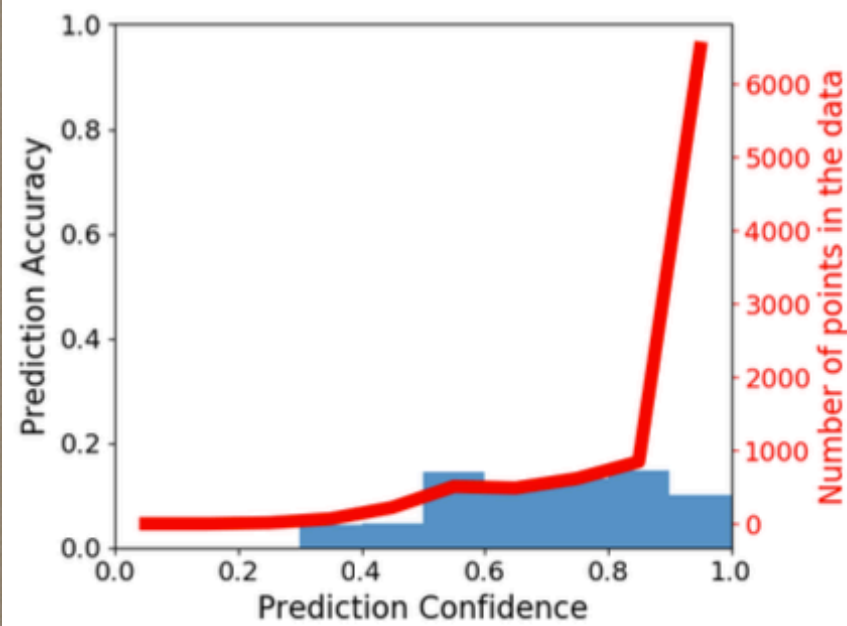
- ◆ Introduction
- ◆ Problem Definition
- ◆ Methodology
- ◆ Experiments
- ◆ Conclusion



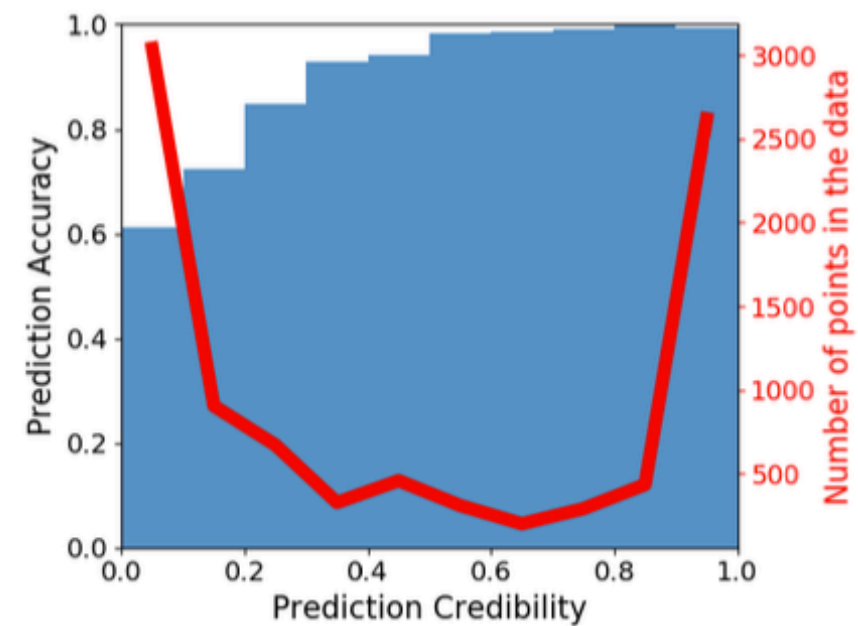
(a) Softmax-MNIST



(b) DkNN-MNIST



(c) Softmax-FGSM



(d) DkNN-FGSM

Reliability Diagrams

DkNN

- ◆ Introduction
- ◆ Problem Definition
- ◆ Methodology
- ◆ Experiments
- ◆ Conclusion

Conclusion

- ◆ Validate the effectiveness of DkNN
- ◆ Hands-on experience on TF2.0

Thank you!