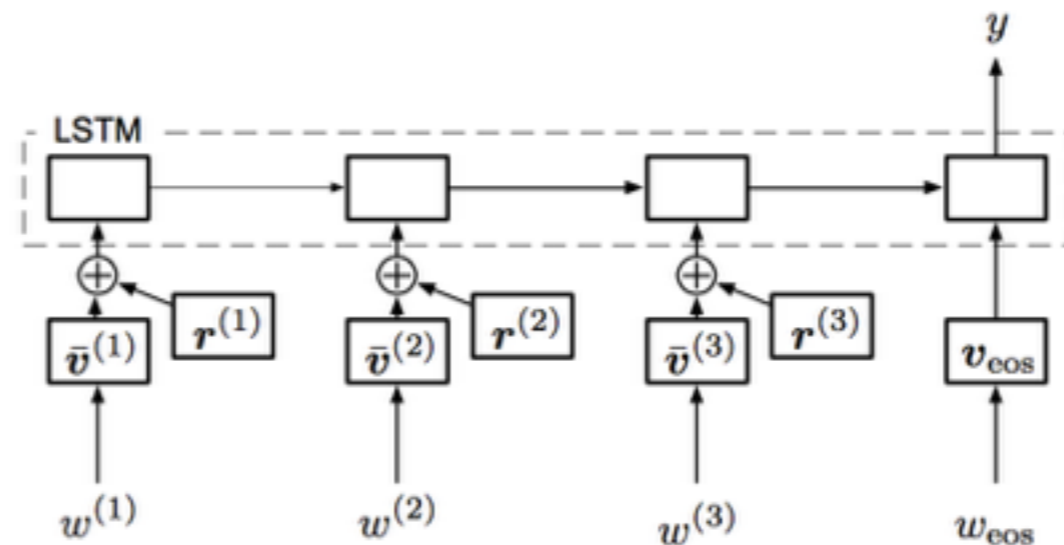


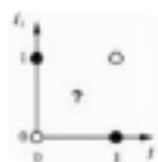
# Virtual Adversarial Training for Semi-Supervised Text Classification

Takeru Miyato, Andrew M. Dai, Ian Goodfellow



Slides By Roei Aharoni  
BIU NLP summer 2016 reading group

# Motivation



**ML Hipster** @ML\_Hipster · 17h

Want to be the best PhD student you can be? Simply make an infinitesimal change to your inputs then take a step in the resulting direction.



# Outline

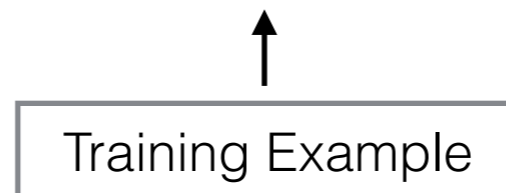
- Introduction to (virtual) adversarial training
- Virtual adversarial training for text classification
- Experimental Setup
- Results (and some analysis)
- Conclusions

# A basic NN training procedure

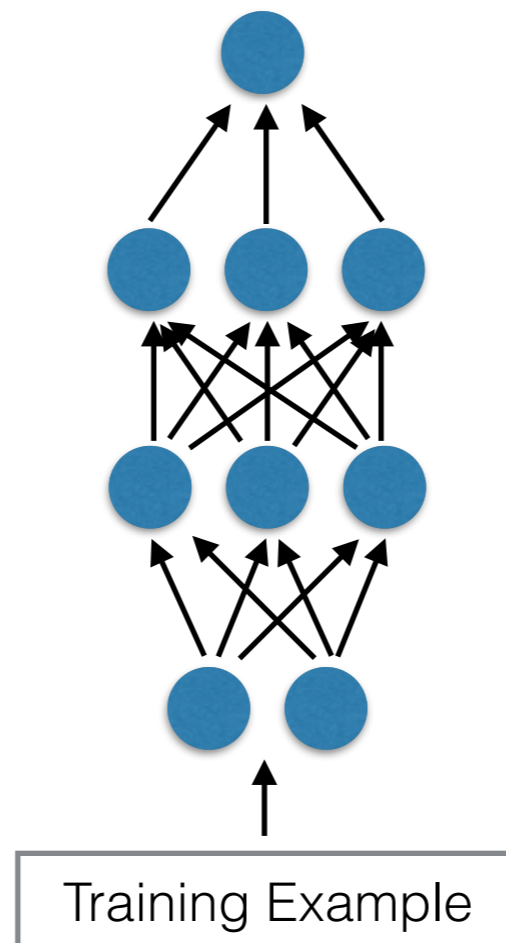
# A basic NN training procedure

Training Example

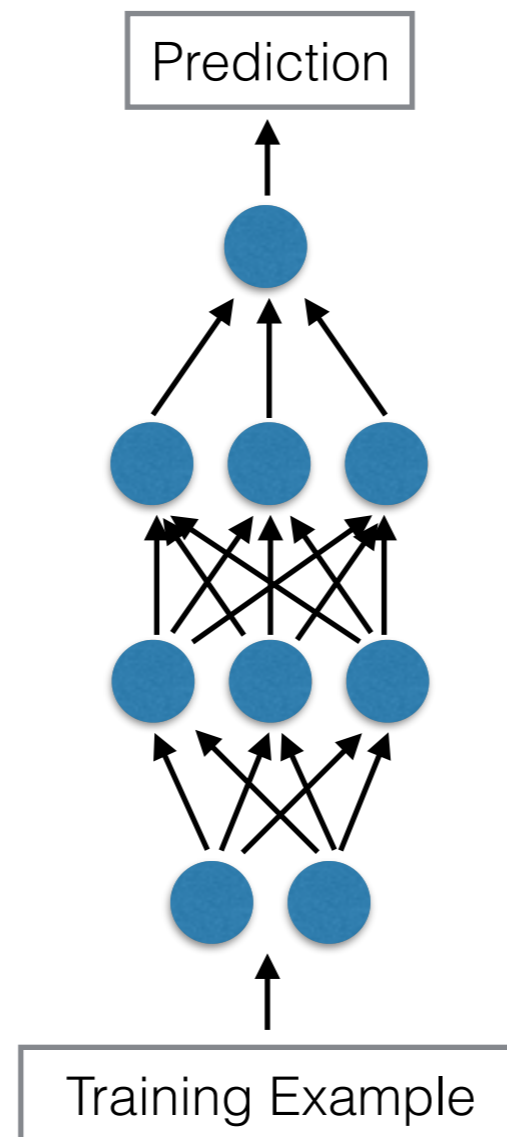
# A basic NN training procedure



# A basic NN training procedure

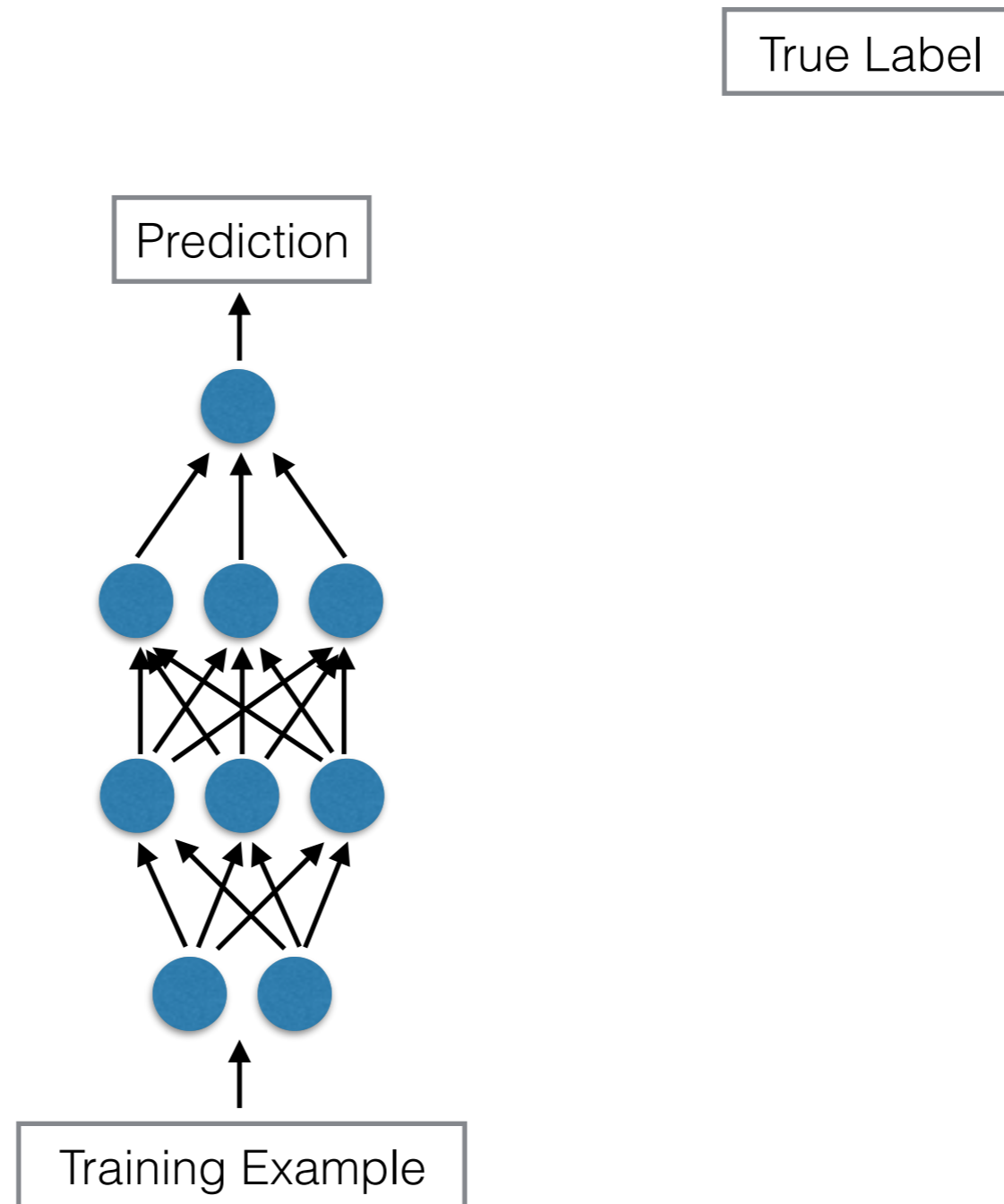


# A basic NN training procedure

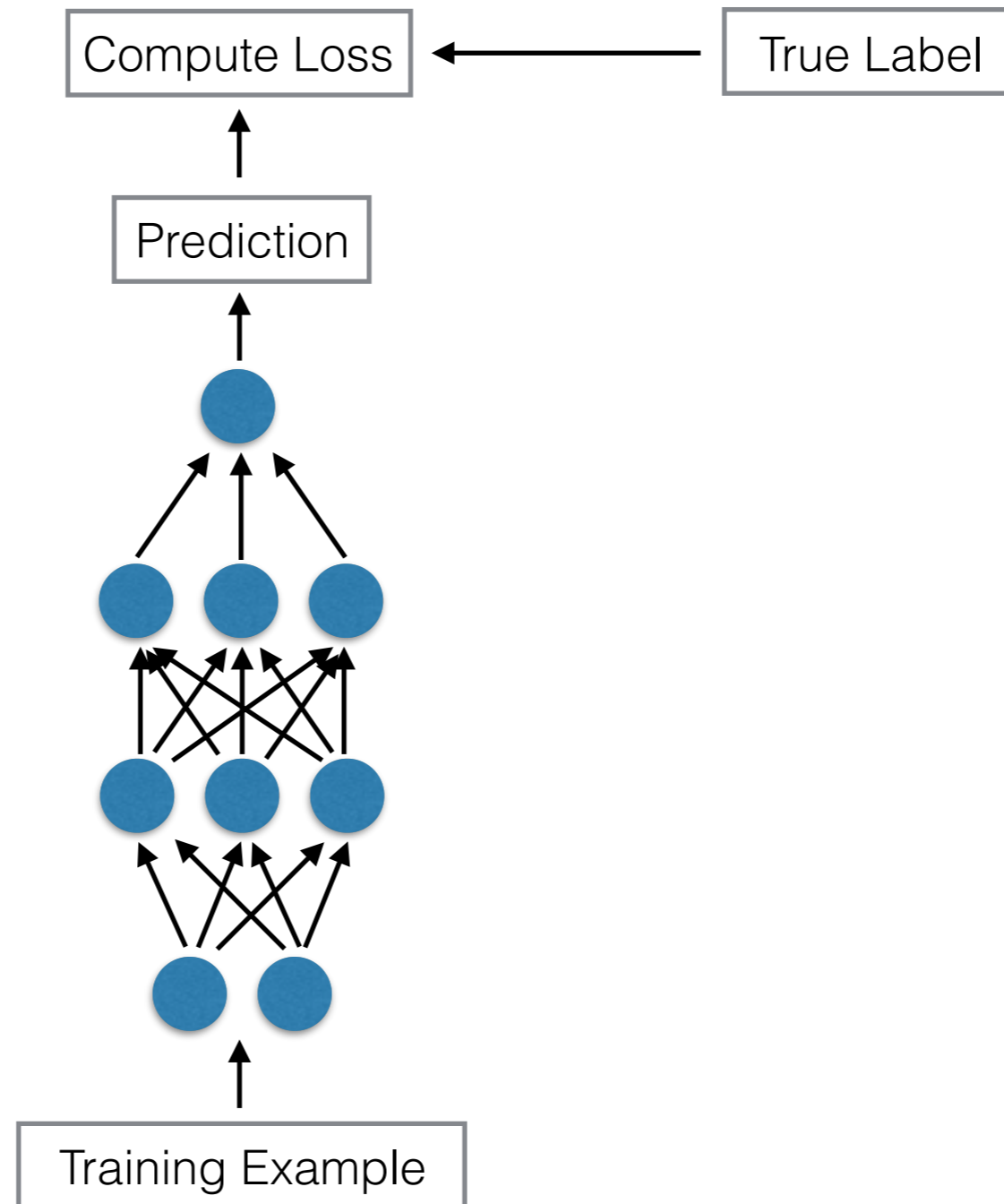




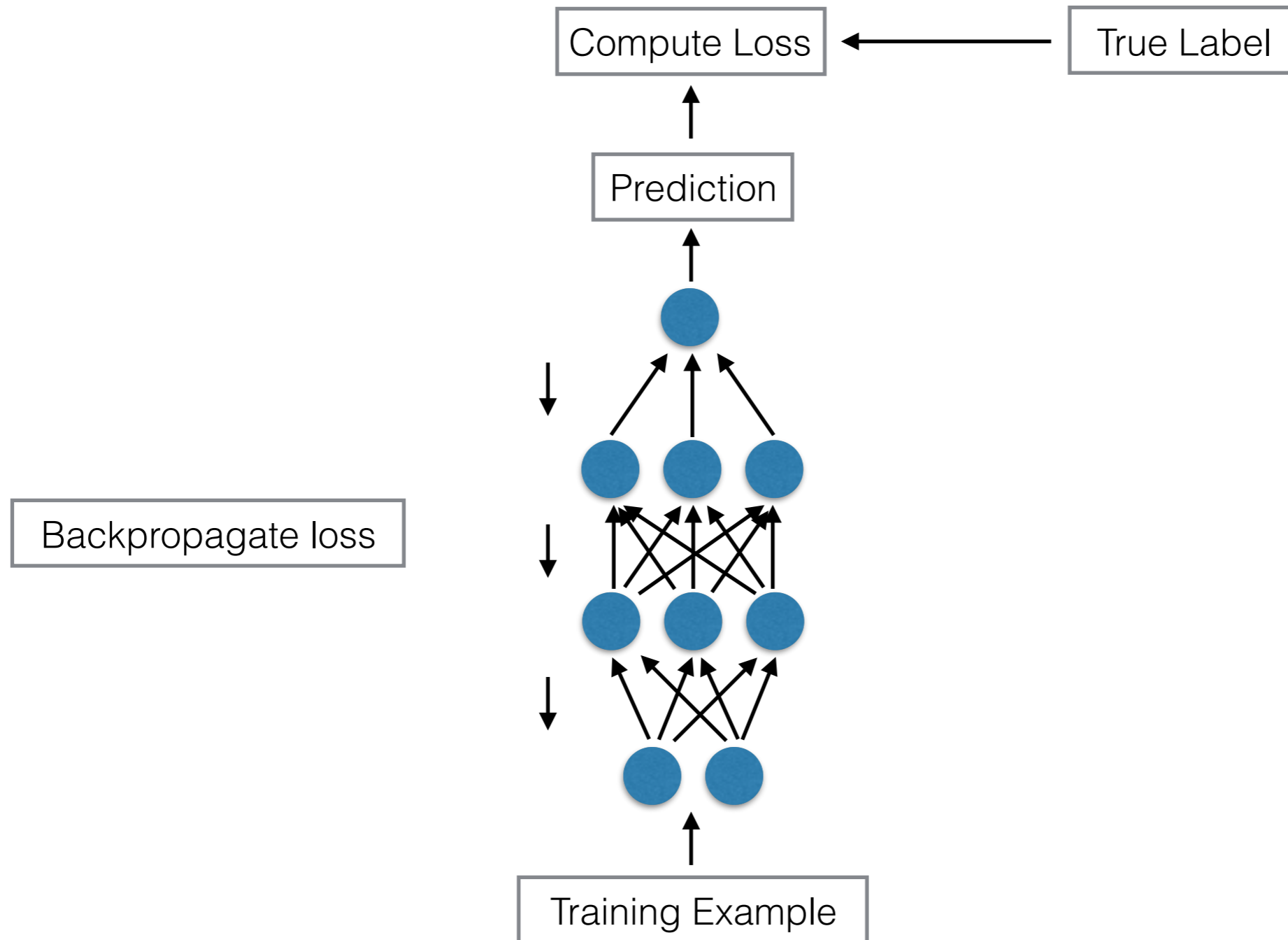
# A basic NN training procedure



# A basic NN training procedure



# A basic NN training procedure



# Adversarial Examples

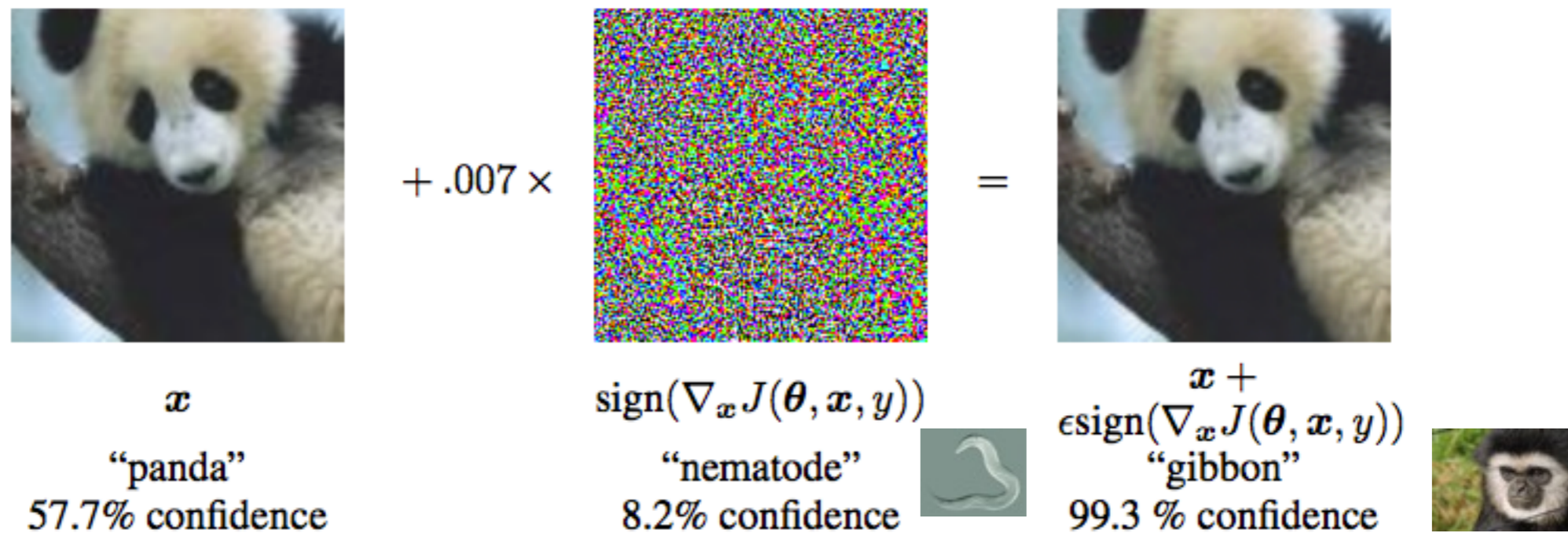


Figure 1: A demonstration of fast adversarial example generation applied to GoogLeNet (Szegedy et al., 2014a) on ImageNet. By adding an imperceptibly small vector whose elements are equal to the sign of the elements of the gradient of the cost function with respect to the input, we can change GoogLeNet’s classification of the image. Here our  $\epsilon$  of .007 corresponds to the magnitude of the smallest bit of an 8 bit image encoding after GoogLeNet’s conversion to real numbers.

# Adversarial Examples

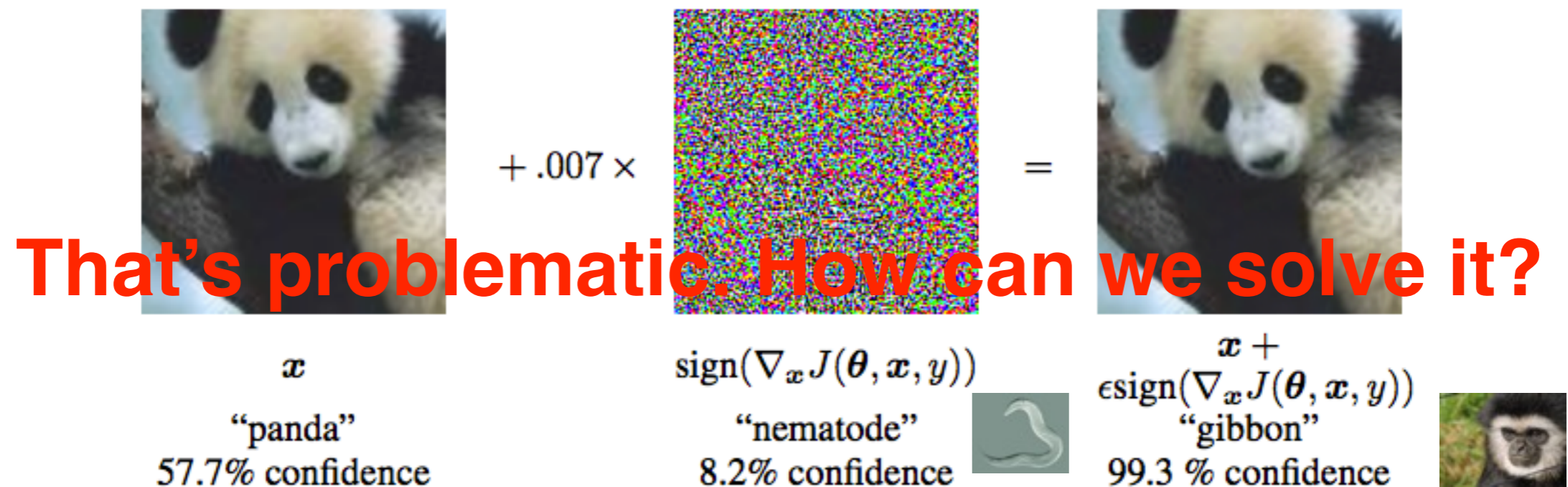


Figure 1: A demonstration of fast adversarial example generation applied to GoogLeNet (Szegedy et al., 2014a) on ImageNet. By adding an imperceptibly small vector whose elements are equal to the sign of the elements of the gradient of the cost function with respect to the input, we can change GoogLeNet's classification of the image. Here our  $\epsilon$  of .007 corresponds to the magnitude of the smallest bit of an 8 bit image encoding after GoogLeNet's conversion to real numbers.

# Adversarial Training for NN's

# Adversarial Training for NN's

Training Example



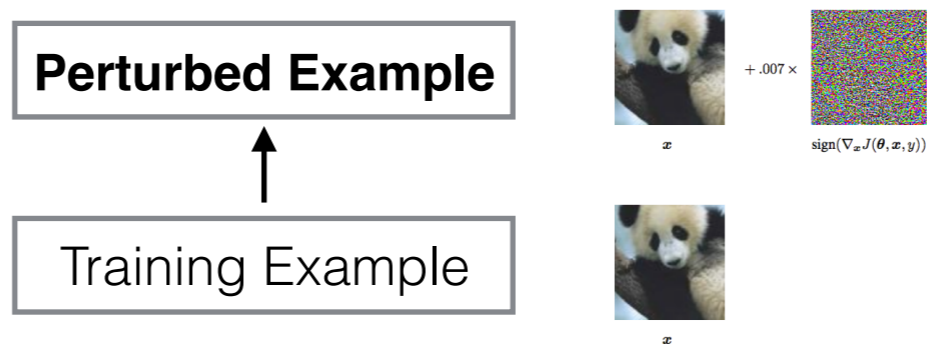
# Adversarial Training for NN's

↑  
Training Example

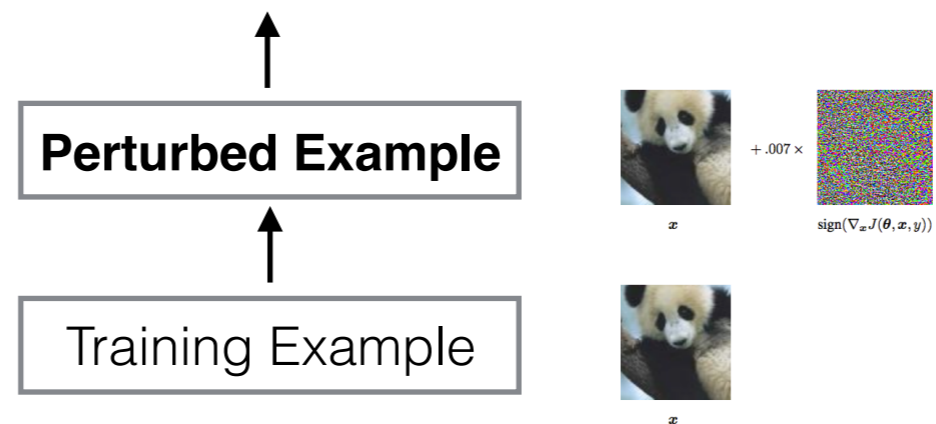




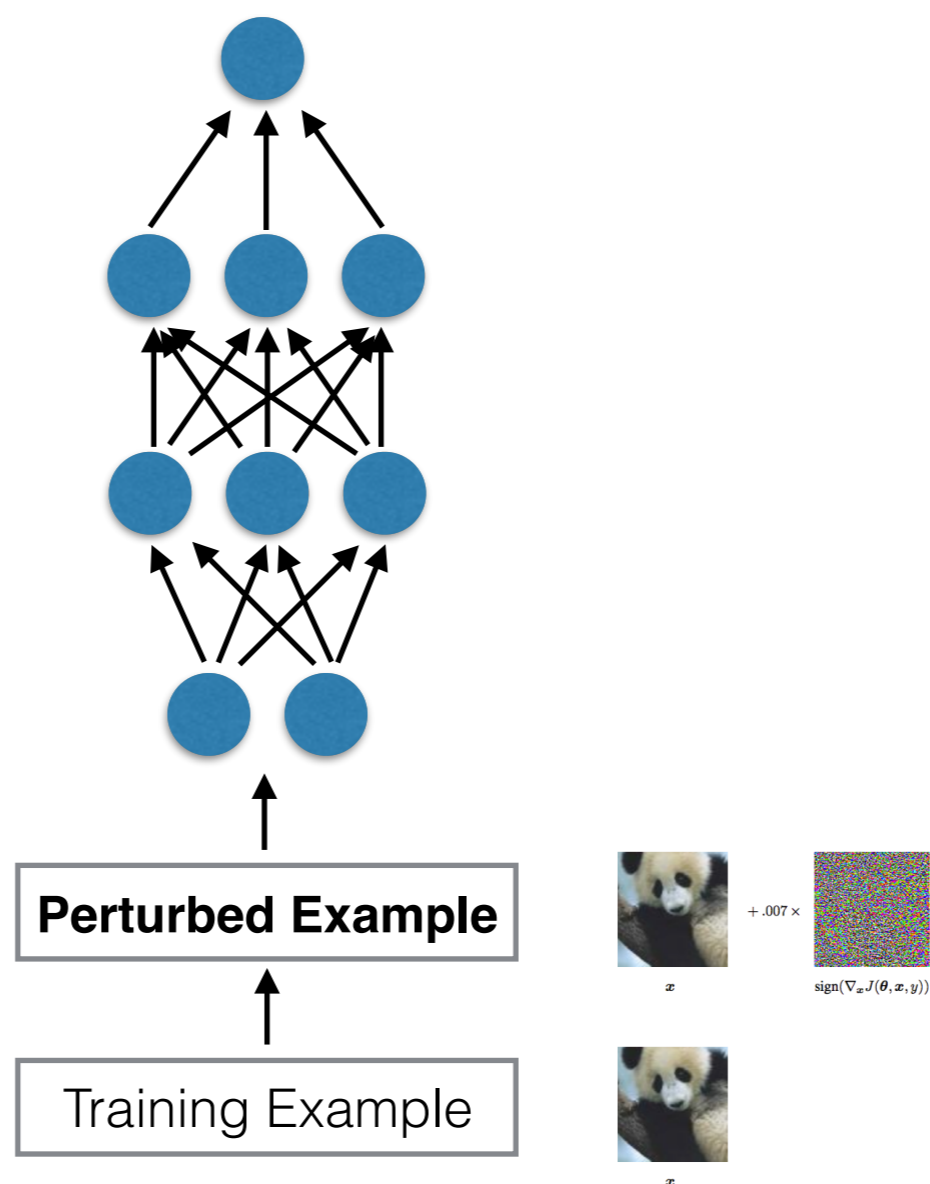
# Adversarial Training for NN's



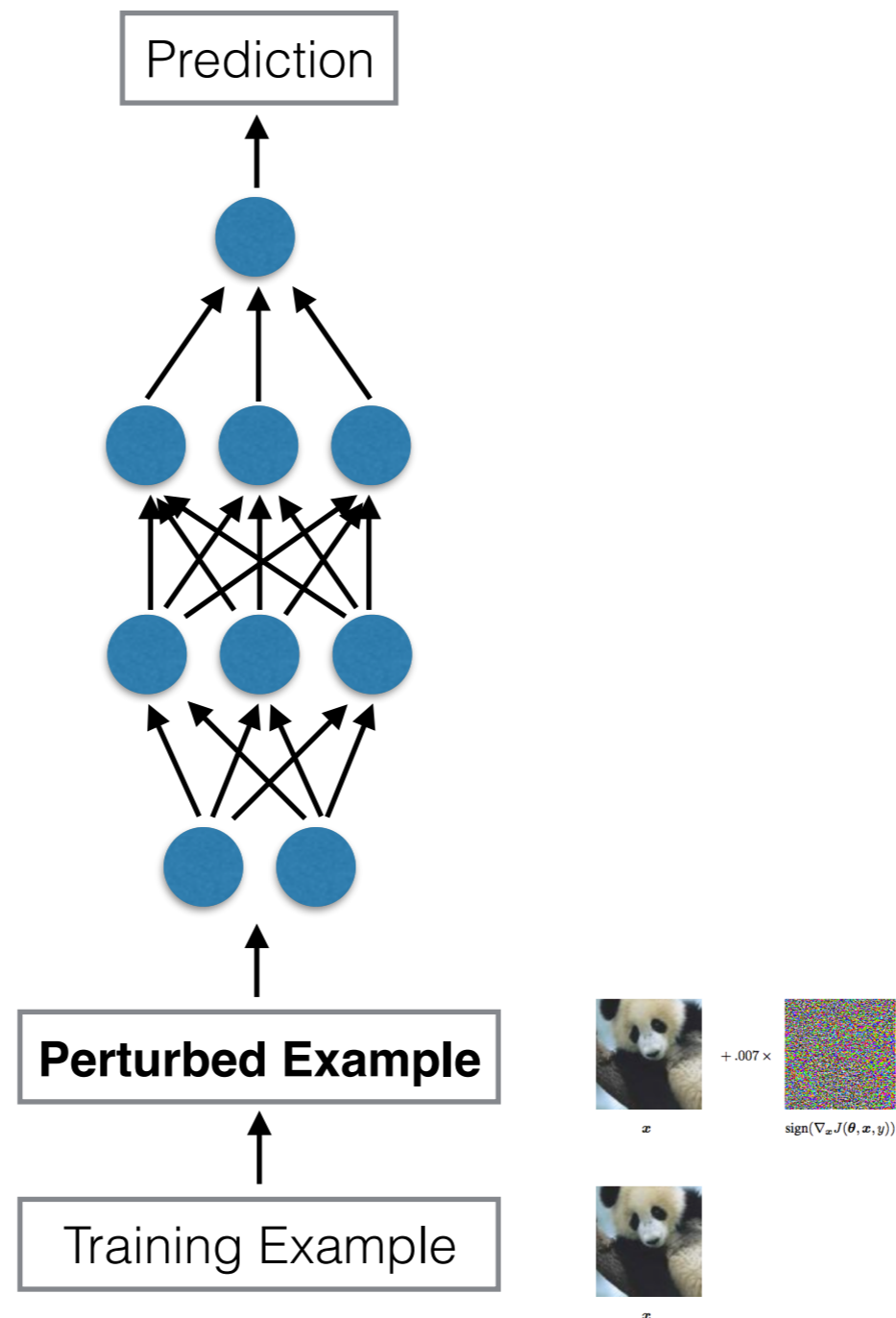
# Adversarial Training for NN's



# Adversarial Training for NN's

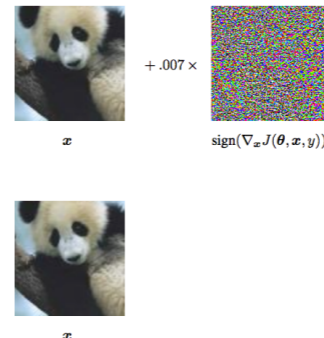
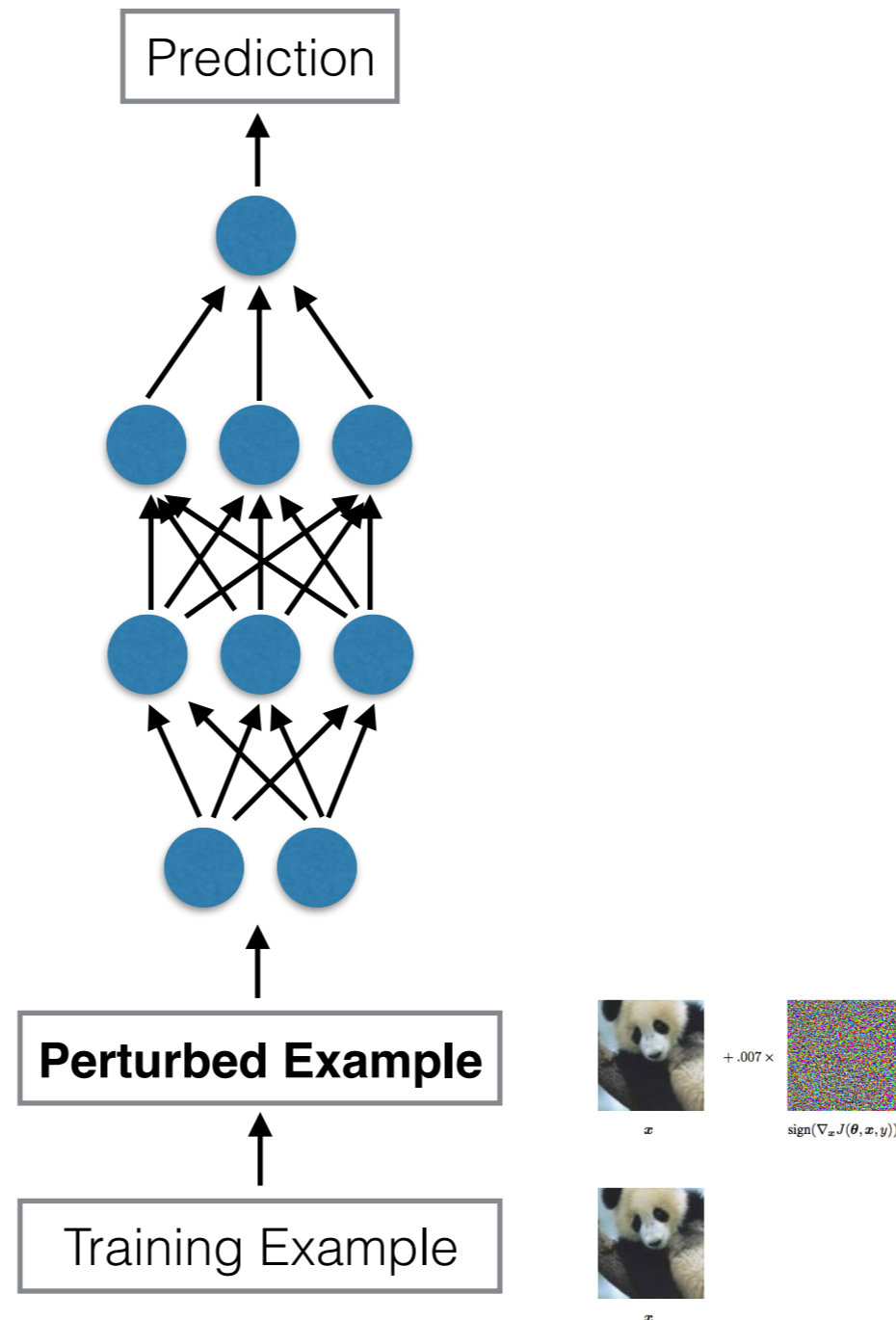


# Adversarial Training for NN's

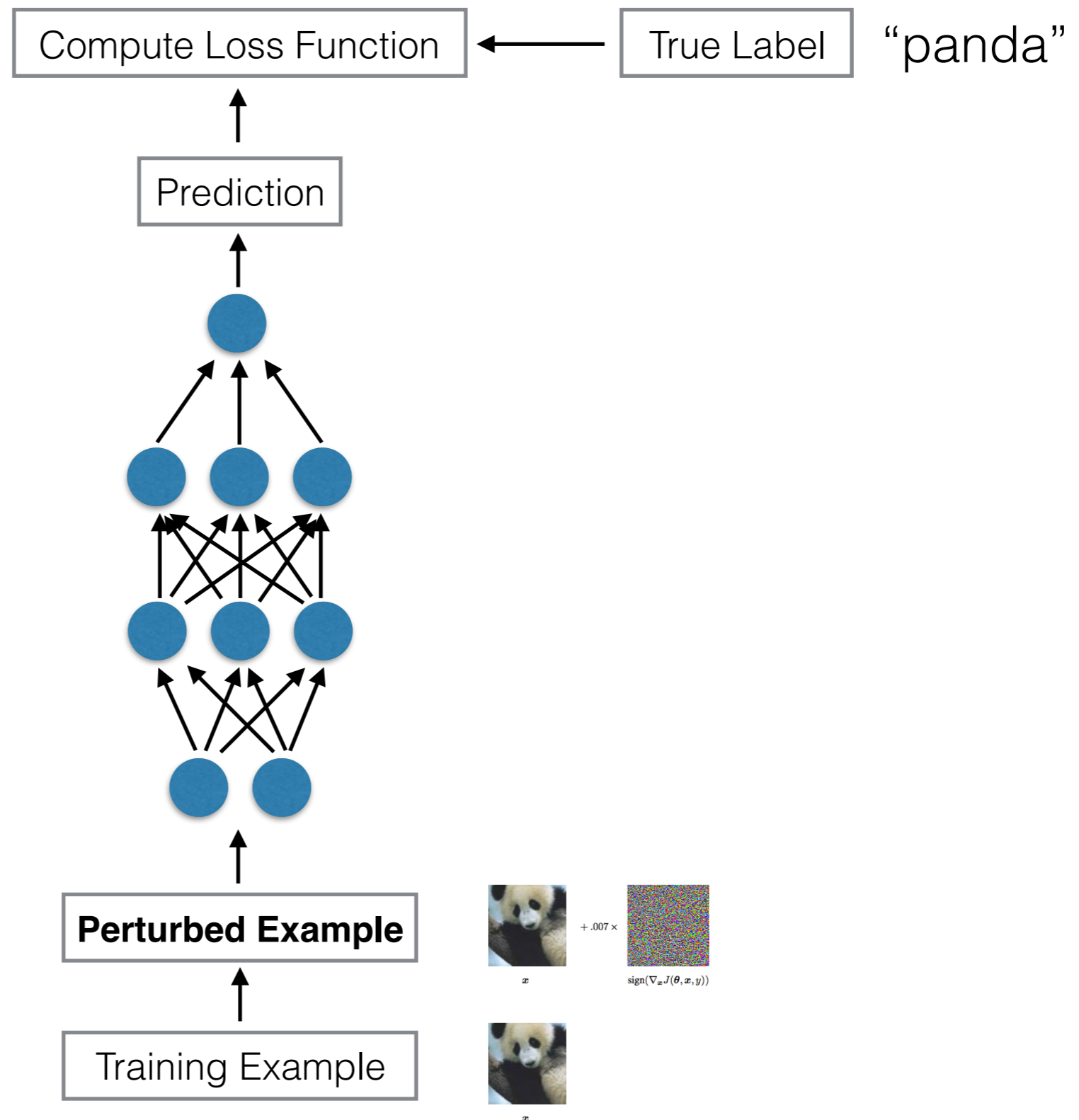


# Adversarial Training for NN's

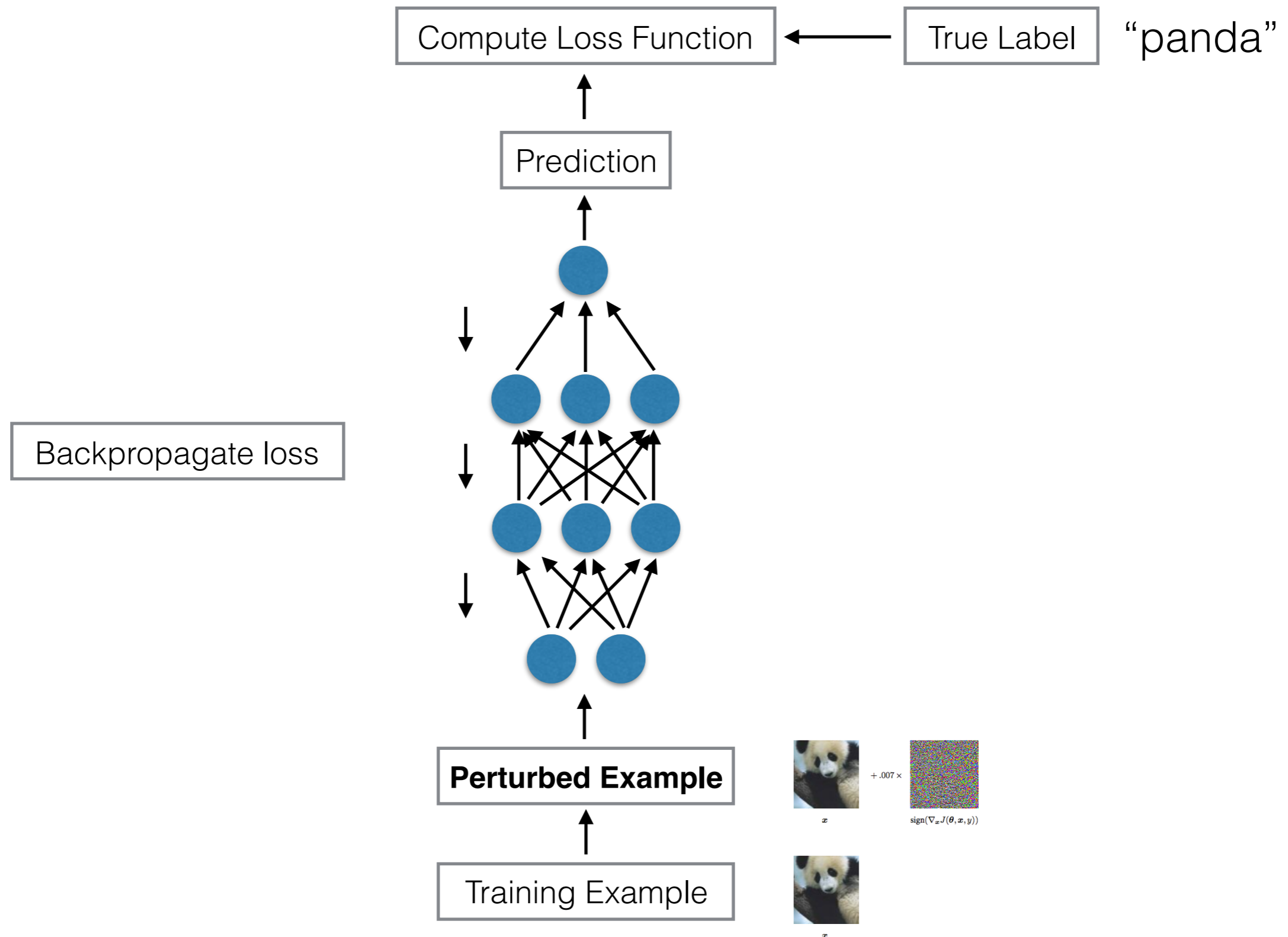
True Label "panda"



# Adversarial Training for NN's



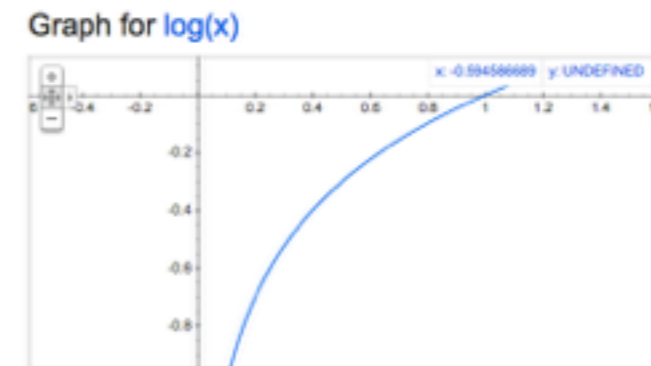
# Adversarial Training for NN's



# Adversarial Training for NN's

- The general addition to the cost function for **adversarial** training:

$$- \min_{\mathbf{r}, \|\mathbf{r}\| \leq \epsilon} \log p(y | \mathbf{x} + \mathbf{r}, \theta)$$



- In practice, take a change depending on the gradient :

$$\mathbf{r}_{\text{adv}} = -\epsilon \mathbf{g} / \|\mathbf{g}\|_2 \text{ where } \mathbf{g} = \nabla_{\mathbf{x}} \log p(y | \mathbf{x}, \theta).$$



# Virtual Adversarial Training

- Extends adversarial training to the semi-supervised regime
- The key idea - make the output **distribution** for an original and perturbed example close to each other
- Enables the use of large amounts of unlabeled data

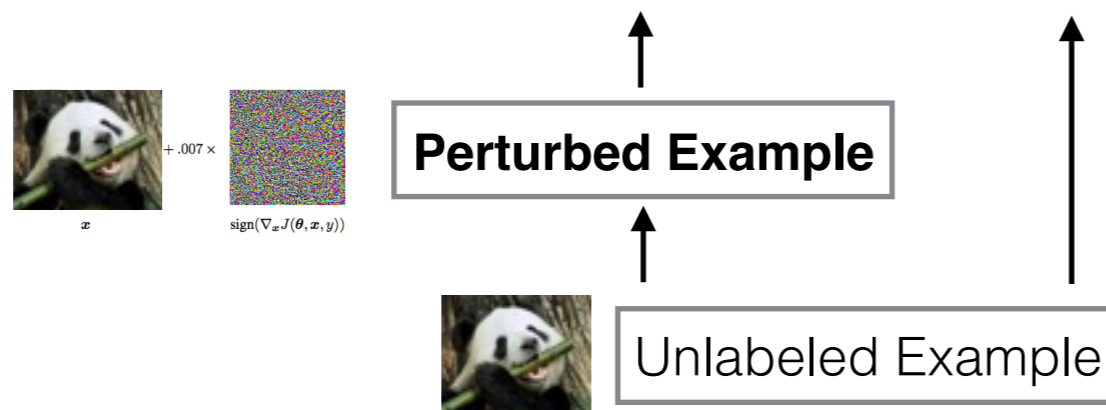
# Virtual Adversarial Training for NN's

# Virtual Adversarial Training for NN's

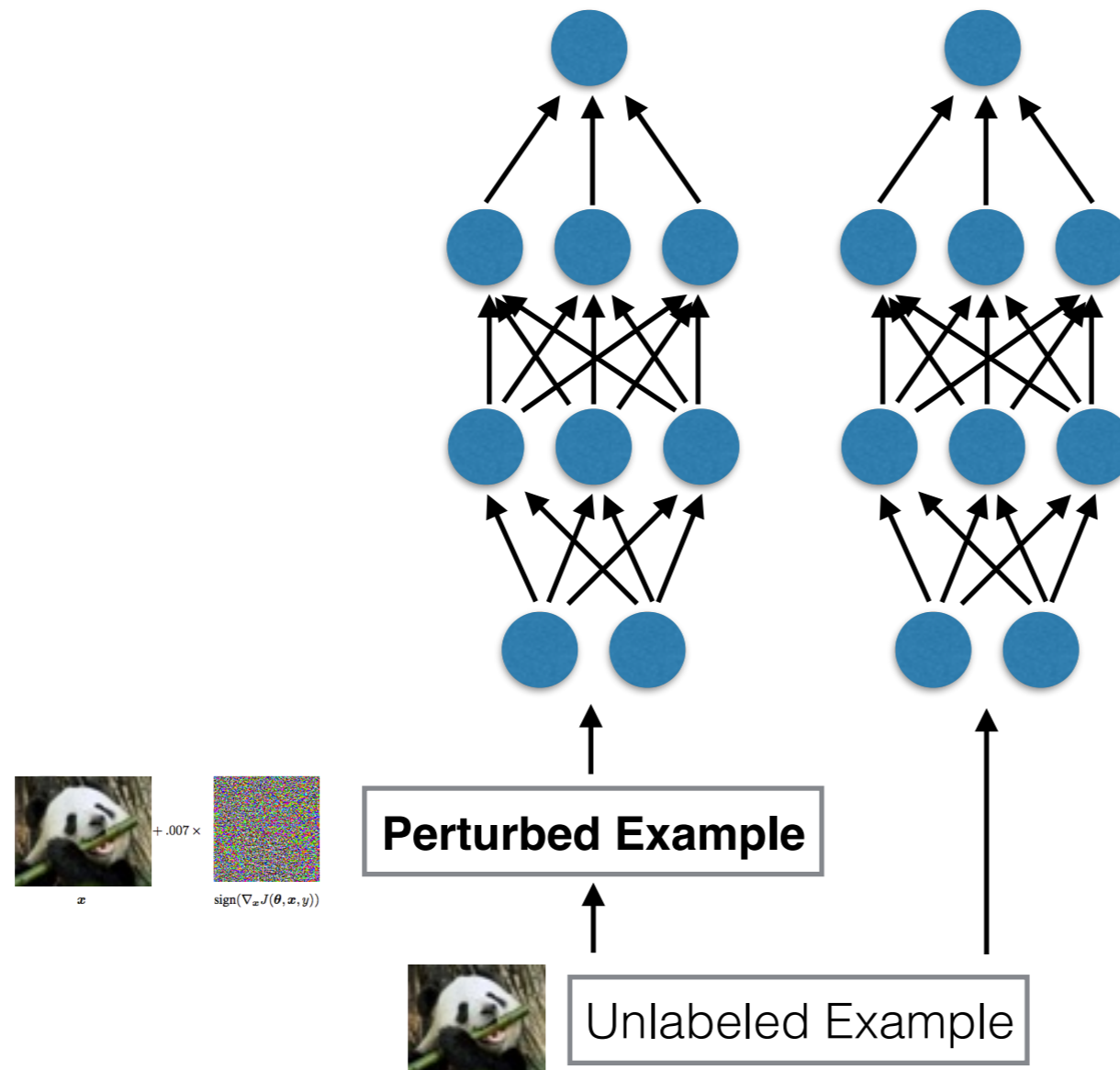


Unlabeled Example

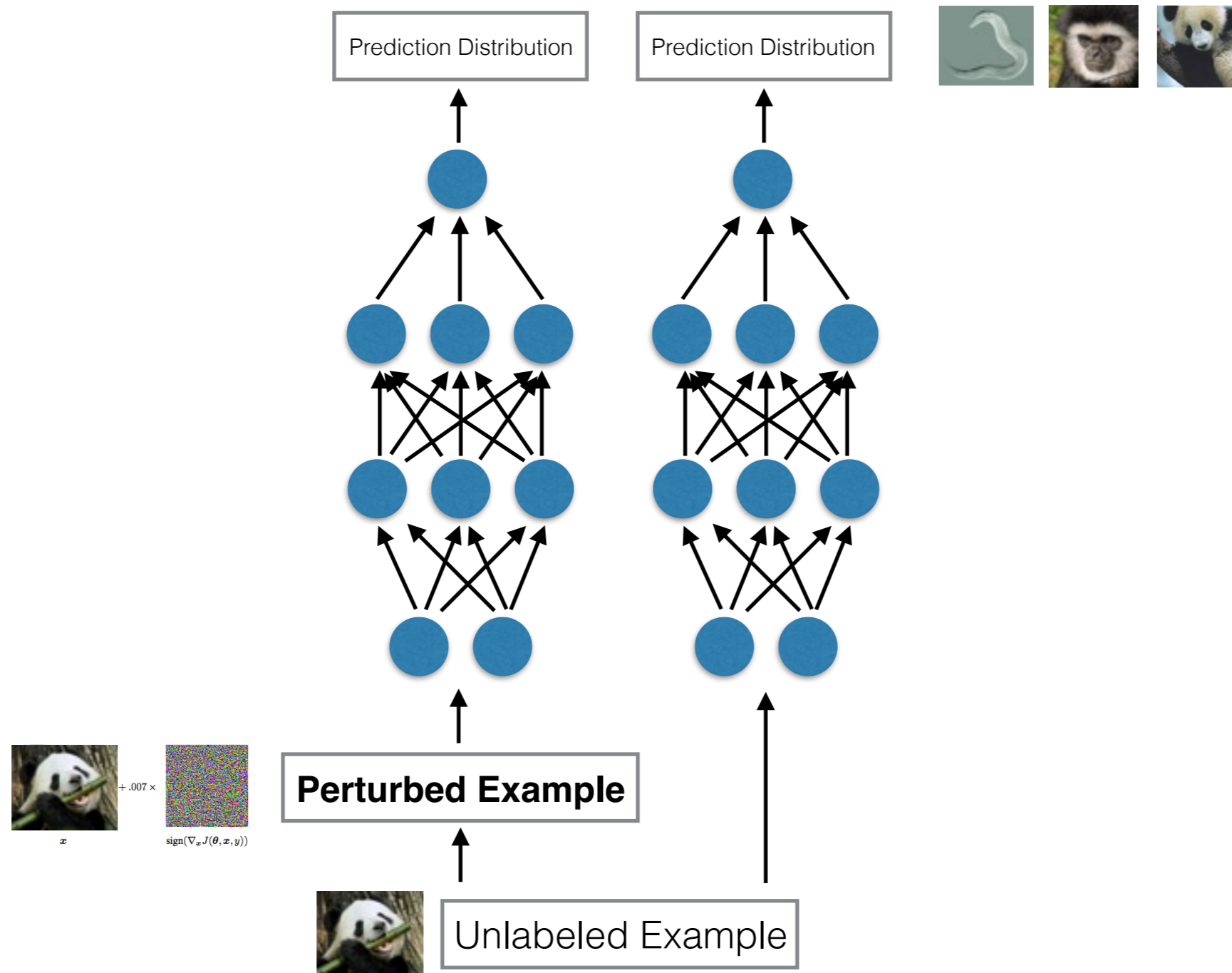
# Virtual Adversarial Training for NN's



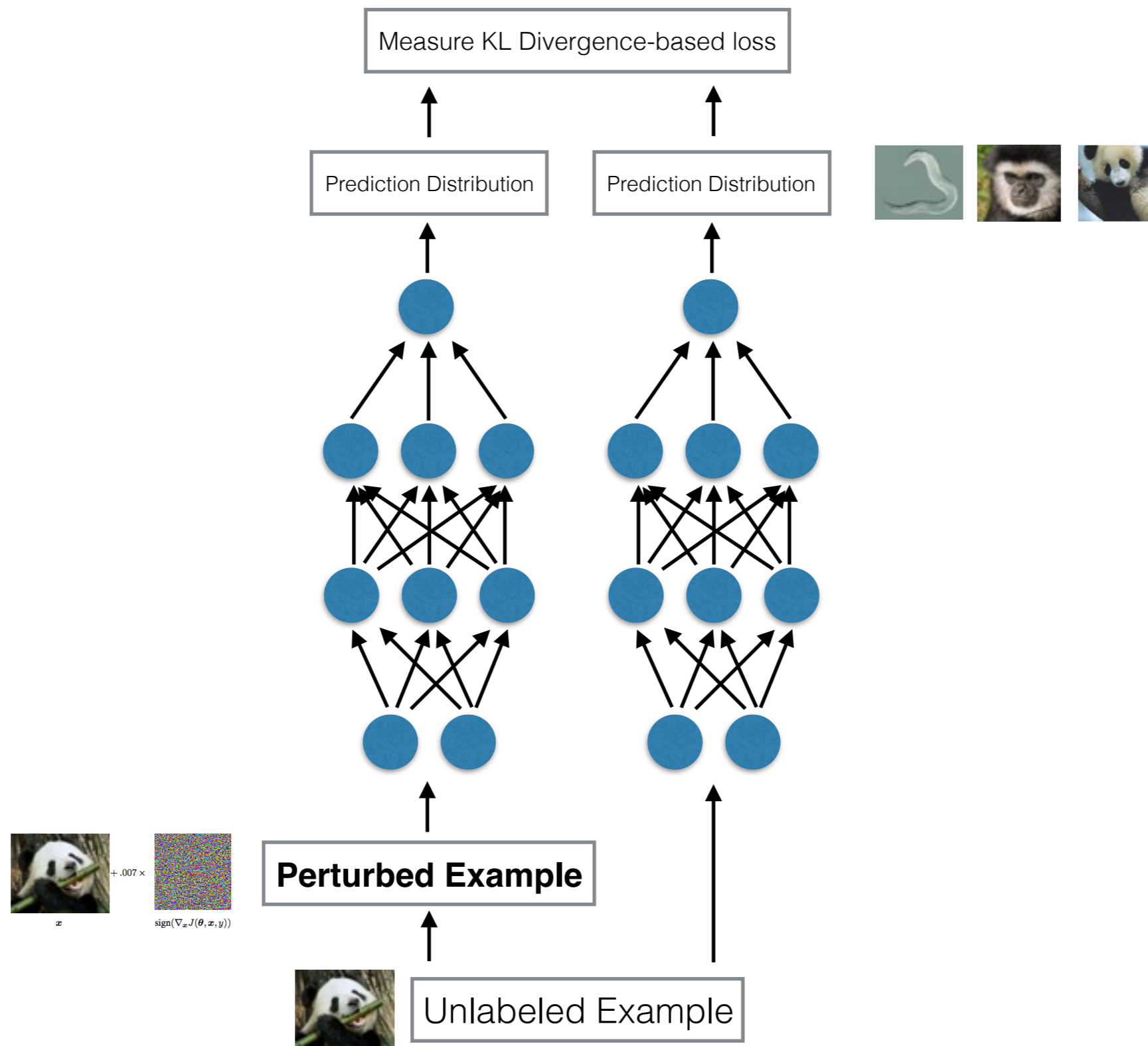
# Virtual Adversarial Training for NN's



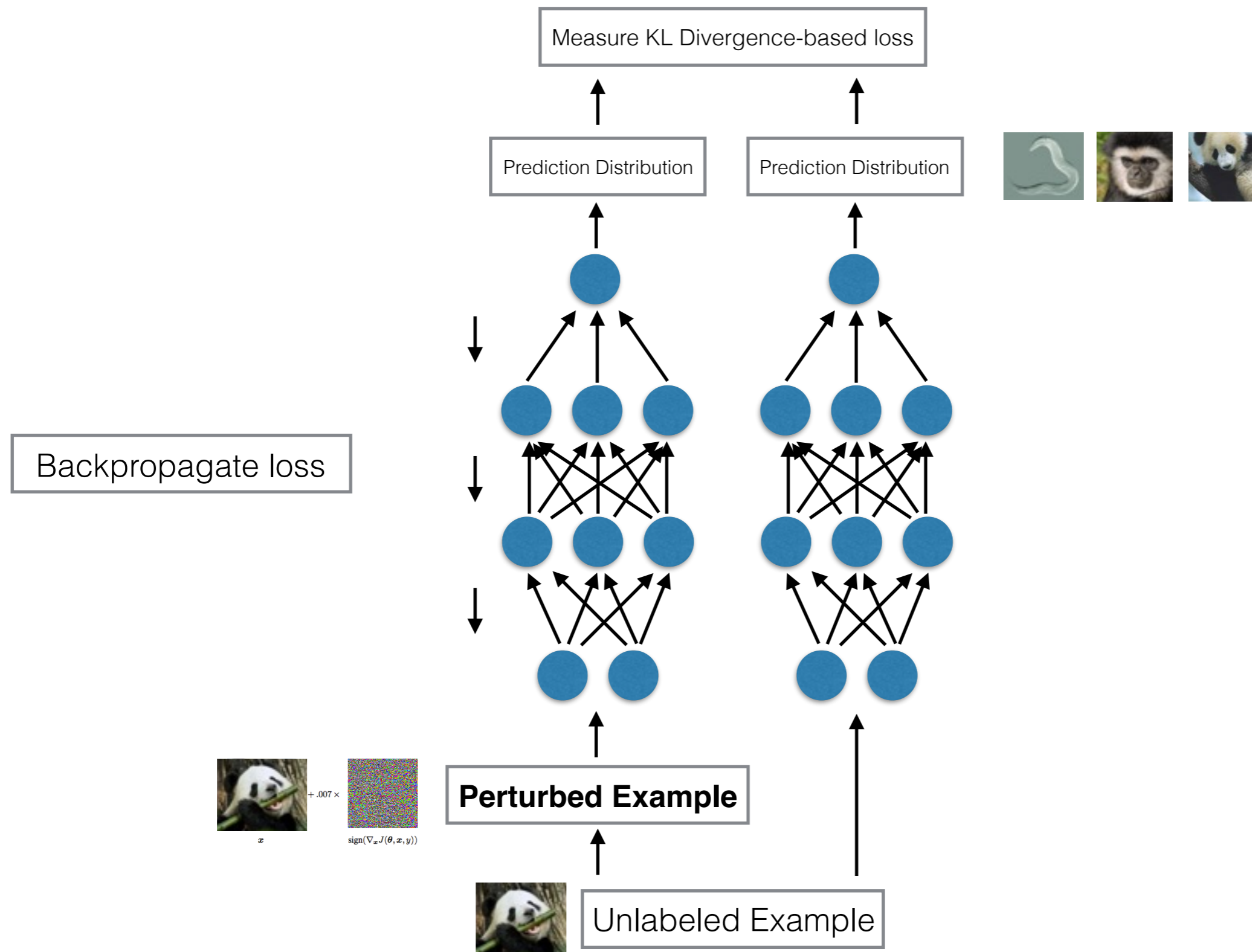
# Virtual Adversarial Training for NN's



# Virtual Adversarial Training for NN's



# Virtual Adversarial Training for NN's





# Virtual Adversarial Training for NN's

- The general addition to the cost function for **virtual adversarial** training:

$$\max_{\mathbf{r}, \|\mathbf{r}\| \leq \epsilon} \text{KL}[p(\cdot | \mathbf{x}, \boldsymbol{\theta}) || p(\cdot | \mathbf{x} + \mathbf{r}, \boldsymbol{\theta})]$$

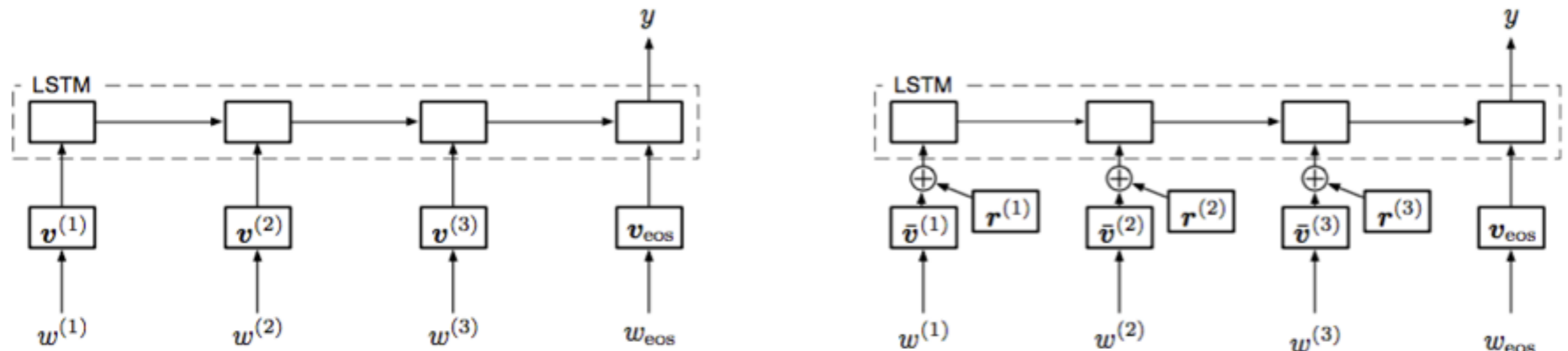
$$D_{\text{KL}}(P||Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

- Again, in practice, there is an efficient way to approximate this (as detailed in Miyato et. al., 2016)

# Model - Adversarial Training for Text Classification

- Adversarial perturbations typically consist of making **small modifications** to very many real-valued inputs (i.e. pixels in the previous examples)
- For text classification, the input is discrete, and usually represented as a series of high-dimensional one-hot vectors (where such small modifications are impossible).
- Solution: define the perturbation on continuous word embeddings instead of discrete word inputs.

# Model - Adversarial Training for Text Classification



(a) LSTM-based text classification model.

(b) The model with perturbed embeddings.

- The perturbation is introduced to normalized embeddings to avoid the network from learning to ignore them:

$$\bar{v}_k = \frac{\mathbf{v}_k - \mathbf{E}(\mathbf{v})}{\sqrt{\text{Var}(\mathbf{v})}} \text{ where } \mathbf{E}(\mathbf{v}) = \sum_{j=1}^K f_j \mathbf{v}_j, \text{ Var}(\mathbf{v}) = \sum_{j=1}^K f_j (\mathbf{v}_j - \mathbf{E}(\mathbf{v}))^2$$

where  $f_i$  is the frequency of the  $i$ -th word, calculated within all training examples.

# Adversarial Training for Text Classification

- As we model the input text as:

$$\mathbf{s} = [\bar{\mathbf{v}}^{(1)}, \bar{\mathbf{v}}^{(2)}, \dots, \bar{\mathbf{v}}^{(T)}]$$

- The perturbation is defined as:

$$\mathbf{r}_{\text{adv}} = -\epsilon \mathbf{g} / \|\mathbf{g}\|_2 \text{ where } \mathbf{g} = \nabla_{\mathbf{s}} \log p(y | \mathbf{s}, \boldsymbol{\theta})$$

- And the addition to the loss function is:

$$L_{\text{adv}}(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{n=1}^N \log p(y_n | \mathbf{s}_n + \mathbf{r}_{\text{adv},n}, \boldsymbol{\theta})$$

# Virtual Adversarial Training for Text Classification

- Here, the perturbation is defined as:

$$\mathbf{r}_{\text{v-adv}} = \epsilon \mathbf{g} / \|\mathbf{g}\|_2 \text{ where } \mathbf{g} = \nabla_{\mathbf{s}+\mathbf{d}} \text{KL} [p(\cdot | \mathbf{s}, \boldsymbol{\theta}) || p(\cdot | \mathbf{s} + \mathbf{d}, \boldsymbol{\theta})]$$

- And the addition to the loss function is then:

$$L_{\text{v-adv}}(\boldsymbol{\theta}) = \frac{1}{N'} \sum_{n'=1}^{N'} \text{KL} [p(\cdot | \mathbf{s}_{n'}, \boldsymbol{\theta}) || p(\cdot | \mathbf{s}_{n'} + \mathbf{r}_{\text{v-adv},n'}, \boldsymbol{\theta})]$$

# Experimental Settings

- 5 datasets:
  - Sentiment classification (binary): IMDB, Rotten Tomatoes, Elec
  - Topic classification (multiclass): DBpedia, RCV1

Table 1: Summary of datasets. Note that unlabeled examples for the Rotten Tomatoes dataset are not provided so we instead use the unlabeled Amazon reviews dataset.

	Classes	Train	Test	Unlabeled	Avg. $T$	Max $T$
IMDB [17]	2	25,000	25,000	50,000	239	2,506
Elec [9]	2	24,792	24,897	197,025	110	5,123
Rotten Tomatoes [23]	2	9596	1066	7,911,684	20	54
DBpedia [14]	14	560,000	70,000	–	49	953
RCV1 [15]	55	15,564	49,838	668,640	153	9,852

# Experimental Settings - Preprocessing

- Treat punctuation as spaces
- Convert words to lower case
- Remove words which appear in only one document
- RCV1 - remove stop words

# Pre-Training Tricks and Hyperparams

- Initialize word embeddings and LSTM weights with RNNLM on labeled and unlabeled examples
- Single layer LSTM, 1024 units (512 for BiLSTM)
- Embedding size: 256(IMDB, BiLSTM)/512(Rest)
- Sampled softmax loss with 1024 candidate samples (?)
- Adam optimization, 256 samples per batch
- 0.5 dropout rate on the word embeddings

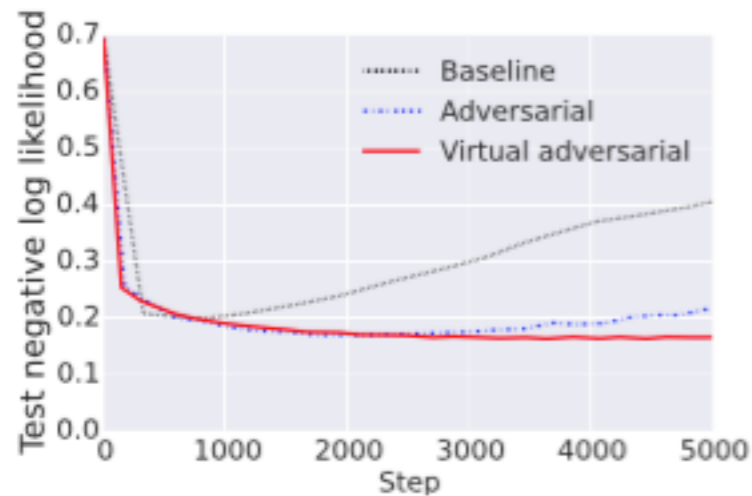


# Classification Model Tricks and Hyperparams

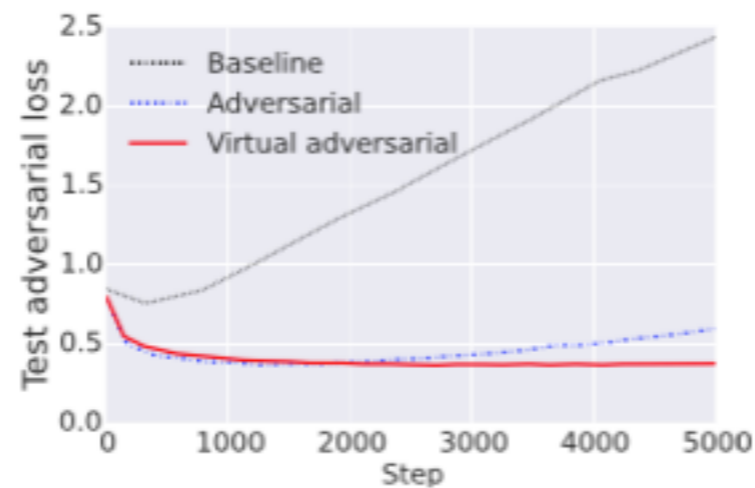
- 1 Hidden layer before softmax, 30(IMDB, Elec, Rotten)/128(Rest) units
- ReLU activation function
- batch size - 64(IMDB, Elec, RCV1)/128(Rest)
- 10k-20k training steps for each model
- Truncated back propagation - stop back propagating after 400 steps
- Generate perturbation after dropout
- Optimize epsilon, dropout rate on validation set

# Results - IMDBB

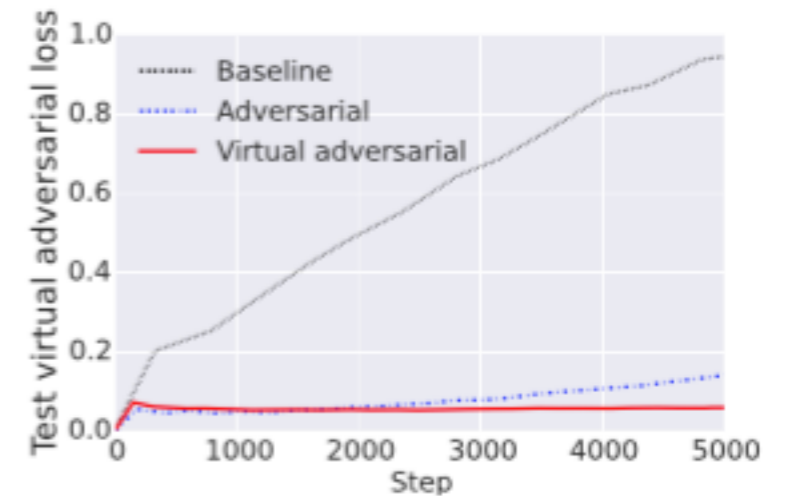
- Adversarial and virtual adversarial training show lower negative log-likelihood
- Virtual adversarial training also improves the adversarial training loss



(a) Negative log likelihood



(b)  $L_{adv}(\theta)$



(c)  $L_{v-adv}(\theta)$

# Results - IMDB

Table 2: Test performance on the IMDB sentiment classification task.

Method	Test error rate
Baseline (without embedding normalization)	7.33%
Baseline	7.39%
Random perturbation with labeled examples	7.20%
Random perturbation with labeled and unlabeled examples	6.78%
Adversarial	6.21%
Virtual Adversarial	<b>5.91%</b>
Adversarial + Virtual Adversarial	6.09%
Virtual Adversarial (on bidirectional LSTM)	<b>5.91%</b>
Adversarial + Virtual Adversarial (on bidirectional LSTM)	6.02%
Full+Unlabeled+BoW [17]	11.11%
Paragraph Vectors [13]	7.42%
SA-LSTM [4]	7.24%
One-hot bi-LSTM (with pretrained embeddings of CNN and bi-LSTM) [10]	5.94%

# Embedding-Based Analysis

Table 3: 10 top nearest neighbors to ‘good’ and ‘bad’ with the word embeddings trained on each method. We used cosine distance for the metric. ‘Baseline’ means training with embedding dropout and ‘Random’ means training with random perturbation with labeled examples. ‘Adversarial’ and ‘Virtual Adversarial’ mean adversarial training and virtual adversarial training.

	‘good’				‘bad’			
	Baseline	Random	Adversarial	Virtual Adversarial	Baseline	Random	Adversarial	Virtual Adversarial
1	great	great	decent	decent	terrible	terrible	terrible	terrible
2	decent	decent	great	great	awful	awful	awful	awful
3	× <u>bad</u>	excellent	nice	nice	horrible	horrible	horrible	horrible
4	excellent	nice	fine	fine	× <u>good</u>	× <u>good</u>	poor	poor
5	Good	Good	entertaining	entertaining	Bad	poor	BAD	BAD
6	fine	× <u>bad</u>	interesting	interesting	BAD	BAD	stupid	stupid
7	nice	fine	Good	Good	poor	Bad	Bad	Bad
8	interesting	interesting	excellent	cool	stupid	stupid	laughable	laughable
9	solid	entertaining	solid	enjoyable	Horrible	Horrible	lame	lame
10	entertaining	solid	cool	excellent	horrendous	horrendous	Horrible	Horrible



# Results - Topic classification: Elec, RCV1

- Improved SOTA on Elec, RCV1, without using CNN's

Table 4: Test performance on the Elec and RCV1 classification tasks

Method	Test error rate	
	Elec	RCV1
Baseline	6.24%	7.40%
Adversarial	5.61%	7.12%
Virtual Adversarial	5.54%	7.05%
Adversarial + Virtual Adversarial	<b>5.40%</b>	6.97%
Virtual Adversarial (on bidirectional LSTM)	5.55%	6.71%
Adversarial + Virtual Adversarial (on bidirectional LSTM)	5.45%	<b>6.68%</b>
One-hot CNN (with pretrained embeddings of CNN) [9]	6.27%	7.71%
One-hot CNN (with pretrained embeddings of CNN and bi-LSTM) [10]	5.82%	7.20%
One-hot bi-LSTM (with pretrained embeddings of CNN and bi-LSTM) [10]	5.55%	8.52%

# Results - Sentiment analysis: Rotten Tomatoes

- Adversarial+Virtual adv. performs equally to SOTA
- Virtual adversarial is weaker than baseline - could be due to small amount of supervised examples, short sentences

Table 5: Test performance on the Rotten Tomatoes sentiment classification task

Method	Test error rate
Baseline	17.9%
Adversarial	16.8%
Virtual Adversarial	19.1%
Adversarial + Virtual Adversarial	<b>16.6%</b>
NBSVM-bigrams[28]	20.6%
CNN (with pretrained embeddings from word2vec Google News)[11]	18.5%
AdaSent (with pretrained embeddings from word2vec Google News)[31]	16.9%
SA-LSTM (with unlabeled data from Amazon reviews)[4]	16.7%

# Results - DBpedia

- Baseline itself improves over SOTA, virtual adversarial performs best

Table 6: Test performance on the DBpedia topic classification task

Method	Test error rate
Baseline (without embedding normalization)	0.87%
Baseline	0.90%
Random perturbation	0.85%
Adversarial	0.79%
Virtual Adversarial	<b>0.76%</b>
Bag-of-words[4]	3.57%
Large-CNN(character-level) [4]	1.73%
SA-LSTM(word-level)[4]	1.41%
N-grams TFIDF [30]	1.31%
SA-LSTM(character-level)[4]	1.19%

# Related Work

- Dropout/Random Noise
- Generative Models
- Pre-Training as semi-supervised learning



# Conclusion

- Adversarial and virtual adversarial training provides good regularization performance for text classification with RNN's
- Provides SOTA results or on-par results for the examined datasets
- “Improved Quality” of word embeddings