# Recent advances in model compression

## Krzysztof J. Geras

NEW YORK UNIVERSITY

Joint work with **Rich Caruana**, **Gregor Urban**, **Abdel-rahman Mohamed**, **Charles Sutton**, **Shengjie Wang**, **Özlem Aslan**, **Samira Ebrahimi Kahou**, **Matthai Philipose** and **Matthew Richardson**

TFML 2017

# Neural networks

# NEURAL NETWORKS

- They work amazingly well.

# NEURAL NETWORKS

- They work amazingly well.
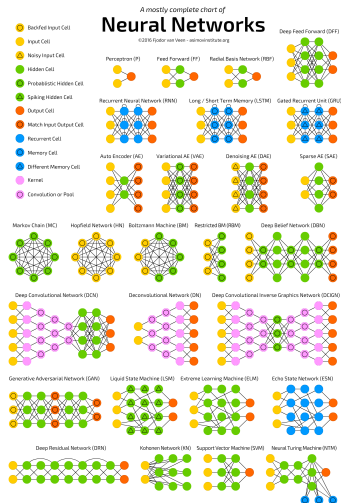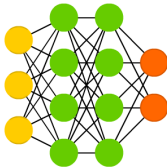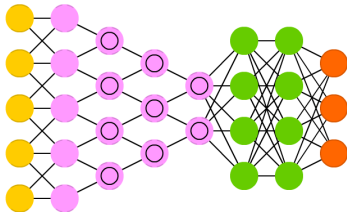- We are largely limited to empirical exploration.

# THE NEURAL NETWORK ZOO



Figure from asimovinstitute.org/neural-network-zoo/.

# Deep Feed Forward (DFF)

# Deep Convolutional Network (DCN)

# Recurrent Neural Network (RNN)
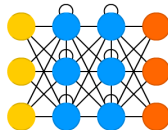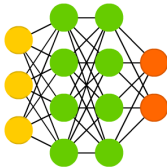
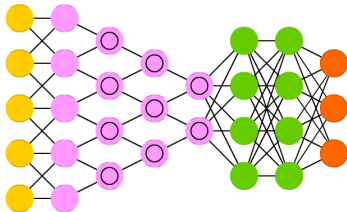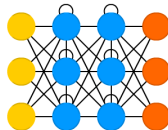Learnability    Representability

Deep Feed Forward (DFF)

Deep Convolutional Network (DCN)

Recurrent Neural Network (RNN)

Learnability ≠ Representability

# MODEL COMPRESSION (AKA KNOWLEDGE DISTILLATION)

- ▶ Idea: take predictions from a big, complex, accurate classifier (a *teacher*) and train a simpler model (a *student*) using them instead of training labels.

# MODEL COMPRESSION (AKA KNOWLEDGE DISTILLATION)

- Idea: take predictions from a big, complex, accurate classifier (a *teacher*) and train a simpler model (a *student*) using them instead of training labels.

- That is, optimise

$$L = -\sum_j \sum_c p(c|\mathbf{x}_j) \log q(c|\mathbf{x}_j),$$

where $p(c|\mathbf{x}_j)$ is teacher's posterior probability of class $c$ given $\mathbf{x}_j$ and $q(c|\mathbf{x}_j)$ is the same for the student.

# MODEL COMPRESSION (AKA KNOWLEDGE DISTILLATION)

▶ Alternatively,

$$L = \lambda \left[ -\sum_j \sum_c p(c|\mathbf{x}_j) \log q(c|\mathbf{x}_j) \right] + (1-\lambda) \left[ -\sum_j \log q(y_j|\mathbf{x}_j) \right],$$

where $p(c|\mathbf{x}_j)$ is teacher's posterior probability of class $c$ given $\mathbf{x}_j$ and $q(c|\mathbf{x}_j)$ is the same for the student.

Why does that work?

Why does that work?

Hypotheses:

Why does that work?

Hypotheses:

- Each example shown to the student model is given with a richer supervision signal.

# MODEL COMPRESSION

Why does that work?

Hypotheses:

- ▶ Each example shown to the student model is given with a richer supervision signal.
- ▶ Cleans noisy labels.

## MODEL COMPRESSION

Why does that work?

Hypotheses:

- ▶ Each example shown to the student model is given with a richer supervision signal.
- ▶ Cleans noisy labels.
- ▶ A way to transfer an inductive bias between models.

# MODEL COMPRESSION

- Large ensemble → single non-convolutional net (Bucila et al., 2006).

# MODEL COMPRESSION

- Large ensemble $\rightarrow$ single non-convolutional net (Bucila et al., 2006).
- Ensemble of deep convolutional nets $\rightarrow$ single shallow non-convolutional net (Ba and Caruana, 2014).

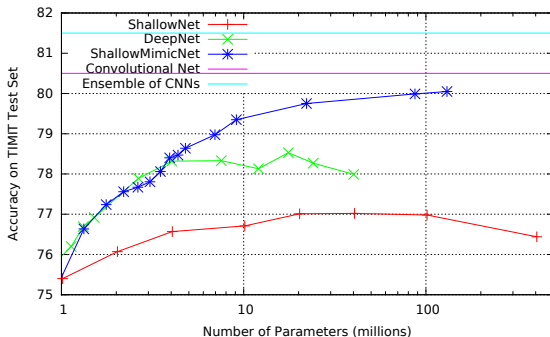# Do deep nets really need to be deep?



Figure from Ba and Caruana (2014).

# MODEL COMPRESSION

- ▶ Large ensemble → single non-convolutional net (Bucila et al., 2006).
- ▶ Ensemble of deep convolutional nets → single shallow non-convolutional net (Ba and Caruana, 2014).
- ▶ Ensemble of deep non-convolutional nets → single deep non-convolutional net (Hinton et al., 2014).

# MODEL COMPRESSION

- Large ensemble $\rightarrow$ single non-convolutional net (Bucila et al., 2006).
- Ensemble of deep convolutional nets $\rightarrow$ single shallow non-convolutional net (Ba and Caruana, 2014).
- Ensemble of deep non-convolutional nets $\rightarrow$ single deep non-convolutional net (Hinton et al., 2014).
- Ensemble of very deep convolutional nets $\rightarrow$ single shallow convolutional net (Urban et al., 2016).

# MODEL COMPRESSION

- Large ensemble → single non-convolutional net (Bucila et al., 2006).
- Ensemble of deep convolutional nets → single shallow non-convolutional net (Ba and Caruana, 2014).
- Ensemble of deep non-convolutional nets → single deep non-convolutional net (Hinton et al., 2014).
- Ensemble of very deep convolutional nets → single shallow convolutional net (Urban et al., 2016).
- Ensemble of deep recurrent nets → single deep convolutional net (Geras et al., 2016).

# DO DEEP CONVOLUTIONAL NETS REALLY NEED TO BE DEEP (OR EVEN CONVOLUTIONAL)?

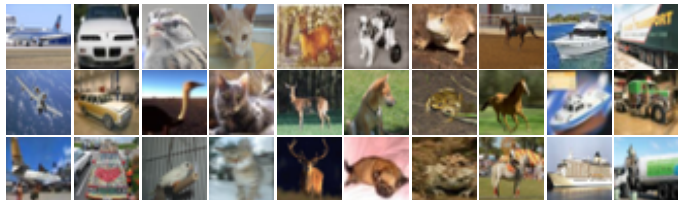- ▶ We know that fully connected nets are compressible.

# DO DEEP CONVOLUTIONAL NETS REALLY NEED TO BE DEEP (OR EVEN CONVOLUTIONAL)?

- We know that fully connected nets are compressible.
- **Question 1.** Can we compress deep convolutional networks into shallow convolutional networks?

# DO DEEP CONVOLUTIONAL NETS REALLY NEED TO BE DEEP (OR EVEN CONVOLUTIONAL)?

- We know that fully connected nets are compressible.
- **Question 1.** Can we compress deep convolutional networks into shallow convolutional networks?
- **Question 2.** Can we compress deep convolutional networks into fully connected networks?

# CIFAR-10 DATA SET



- ▸ Labelled subset of the Tiny 80M images data set.
- ▸ 60k 32x32 RGB images.
- ▸ 10 classes: airplane, automobile, bird, cat, deer, dog, frog, horse, sheep, truck.
- ▸ Each class contains 6k images.

- The teacher: 8 convolutional layers and 2 fully connected layers.

- ► The teacher: 8 convolutional layers and 2 fully connected layers.
- ► Various possible student architectures.

- The teacher: 8 convolutional layers and 2 fully connected layers.
- Various possible student architectures.
- We need to be extremely careful.
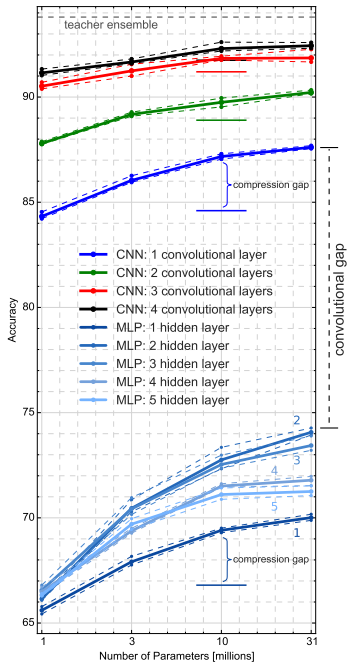
- The teacher: 8 convolutional layers and 2 fully connected layers.
- Various possible student architectures.
- We need to be extremely careful.
- We use Bayesian optimisation to find the best hyperparameters.

Figure legend:
- CNN: 1 convolutional layer
- CNN: 2 convolutional layers
- CNN: 3 convolutional layers
- CNN: 4 convolutional layers
- MLP: 1 hidden layer
- MLP: 2 hidden layer
- MLP: 3 hidden layer
- MLP: 4 hidden layer
- MLP: 5 hidden layer

Axis labels: Accuracy (y-axis), Number of Parameters [millions] (x-axis)

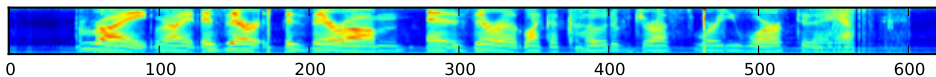Annotations: teacher ensemble, compression gap, convolutional gap
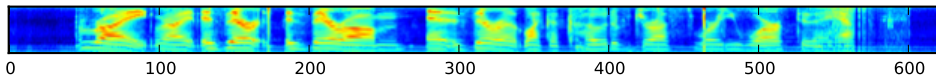
Deep convolutional nets really need to be deep.

Deep convolutional nets really need to be deep. And convolutional.

Deep convolutional nets really need to be deep. And convolutional. But perhaps not that deep.

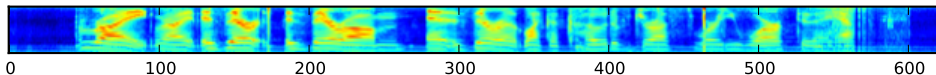# Speech recognition $\overset{?}{\approx}$ object recognition

# SPEECH RECOGNITION $\overset{?}{\approx}$ OBJECT RECOGNITION
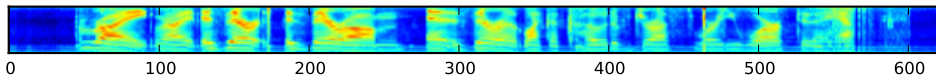


Speech recognition framework:

# SPEECH RECOGNITION $\overset{?}{\approx}$ OBJECT RECOGNITION



Speech recognition framework:

- ▶ Sample many windows of speech. Train a classifier.

Speech recognition framework:

- Sample many windows of speech. Train a classifier.
- Use decoding to get words.

- A benchmark for speech recognition.

# THE SWITCHBOARD DATA SET

- A benchmark for speech recognition.
- Very large, 309 hours of speech, 18 GB.

## THE SWITCHBOARD DATA SET

- A benchmark for speech recognition.
- Very large, 309 hours of speech, 18 GB.
- We sample training examples of size $31 \times 41$, 9000 output classes.

# CNNs for Speech

| |
|---|
| softmax |
| fully connected, 2048 |
| fully connected, 2048 |
| fully connected, 2048 |
| fully connected, 2048 |
| convolution, 7×7, 324 |
| max pooling, 3×1 |
| convolution, 7×7, 324 |
| input (31x41) |

# CNNs for speech

| softmax |
| --- |
| fully connected, 4096 |
| fully connected, 4096 |
| max pooling, 2×2 |
| convolution, 3×3, 384 |
| convolution, 3×3, 384 |
| convolution, 3×3, 384 |
| max pooling, 2×2 |
| convolution, 3×3, 192 |
| convolution, 3×3, 192 |
| convolution, 3×3, 192 |
| max pooling, 2×2 |
| convolution, 3×3, 96 |
| convolution, 3×3, 96 |
| input (31x41) |

| softmax |
| --- |
| fully connected, 2048 |
| fully connected, 2048 |
| fully connected, 2048 |
| fully connected, 2048 |
| convolution, 7×7, 324 |
| max pooling, 3×1 |
| convolution, 7×7, 324 |
| input (31x41) |

# CNNS FOR SPEECH

| softmax |
|---|
| fully connected, 4096 |
| fully connected, 4096 |

| max pooling, 2×2 |
|---|
| convolution, 3×3, 384 |
| convolution, 3×3, 384 |
| convolution, 3×3, 384 |

| max pooling, 2×2 |
|---|
| convolution, 3×3, 192 |
| convolution, 3×3, 192 |
| convolution, 3×3, 192 |

| max pooling, 2×2 |
|---|
| convolution, 3×3, 96 |
| convolution, 3×3, 96 |

| input (31x41) |
|---|

| softmax |
|---|
| fully connected, 2048 |
| fully connected, 2048 |
| fully connected, 2048 |
| fully connected, 2048 |
| convolution, 7×7, 324 |

| max pooling, 3×1 |
|---|
| convolution, 7×7, 324 |

| input (31x41) |
|---|

**Sainath et al.-style CNN**          **vision-style CNN**

# CNNs vs LSTMs on the Switchboard data set

|  | frame error rate | word error rate |
|---|---|---|
| Sainath et al.-style CNN | 37.93% | 15.5 |
| vision-style CNN | 35.51% | 14.1 |
| LSTM | 34.15% | 14.4 |

# CNNS VS LSTMS ON THE SWITCHBOARD DATA SET

| | frame error rate | word error rate |
|---|---|---|
| Sainath et al.-style CNN | 37.93% | 15.5 |
| vision-style CNN | 35.51% | 14.1 |
| LSTM | 34.15% | 14.4 |

# CNNs vs LSTMs on the Switchboard data set

| | frame error rate | word error rate |
|---|---|---|
| Sainath et al.-style CNN | 37.93% | 15.5 |
| vision-style CNN | 35.51% | 14.1 |
| LSTM | 34.15% | 14.4 |

# CNNs vs LSTMs for speech
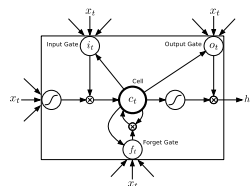


Figure from deeplearning.net



Figure from Graves et al. (2013)

# CAN WE DO BETTER?

- Different network structures $\rightarrow$ different inductive biases.

# CAN WE DO BETTER?

- Different network structures $\rightarrow$ different inductive biases.
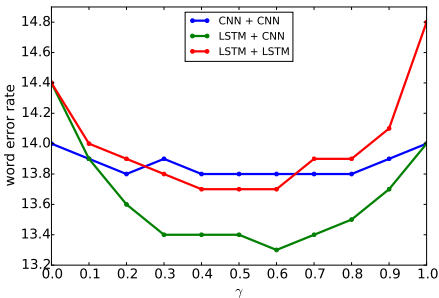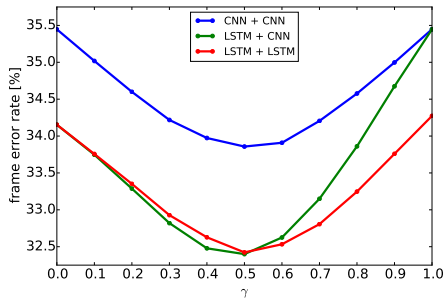- **Can we have two models in one?**

# CAN WE DO BETTER?

- Different network structures $\rightarrow$ different inductive biases.
- **Can we have two models in one?**
- Yes, there is an easy way to do this - ensembling.

$$p(y|\mathbf{x}_i) = \gamma p_{\text{LSTM}}(y|\mathbf{x}_i) + (1 - \gamma)p_{\text{CNN}}(y|\mathbf{x}_i)$$
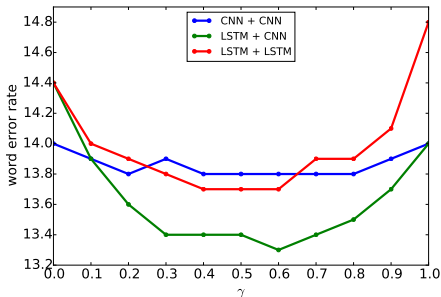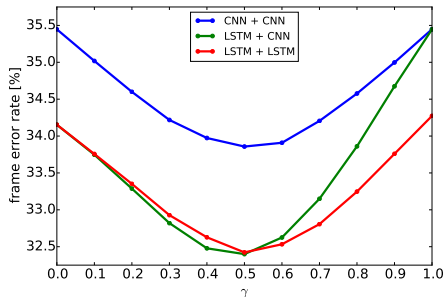
# ENSEMBLING

$$p(y|\mathbf{x}_i) = \gamma p_{\text{LSTM}}(y|\mathbf{x}_i) + (1 - \gamma)p_{\text{CNN}}(y|\mathbf{x}_i)$$

# ENSEMBLING

$$p(y|\mathbf{x}_i) = \gamma p_{\text{LSTM}}(y|\mathbf{x}_i) + (1 - \gamma)p_{\text{CNN}}(y|\mathbf{x}_i)$$



**Big issue:** LSTM is **6** times slower than the CNN. We need to have two models in one CNN.

- Very large data set, 309 hours of speech, 18 GB.

- Very large data set, 309 hours of speech, 18 GB.
- 31×41 inputs, 9000 output classes.

- Very large data set, 309 hours of speech, 18 GB.
- $31 \times 41$ inputs, 9000 output classes. $\rightarrow$ Predictions would take 3.6 TB.
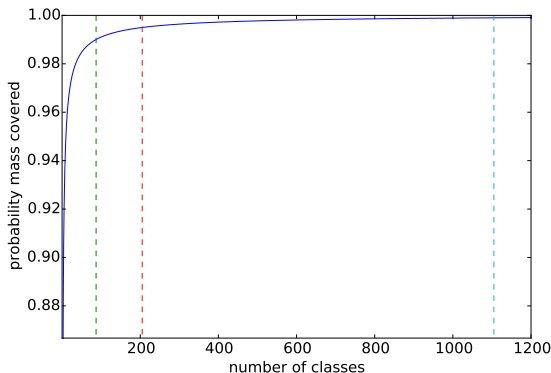
# HOW TO DO COMPRESSION WITH SWITCHBOARD

$$M(C) = \frac{1}{|\{\mathbf{x}_i\}|} \sum_{\mathbf{x}_i} \sum_{y \in \text{TOP}_C(\mathbf{x}_i)} p(y|\mathbf{x}_i).$$

# HOW TO DO COMPRESSION WITH SWITCHBOARD

$$\mathrm{M}(C) = \frac{1}{|\{\mathbf{x}_i\}|} \sum_{\mathbf{x}_i} \sum_{y \in \mathrm{TOP}_C(\mathbf{x}_i)} p(y|\mathbf{x}_i).$$
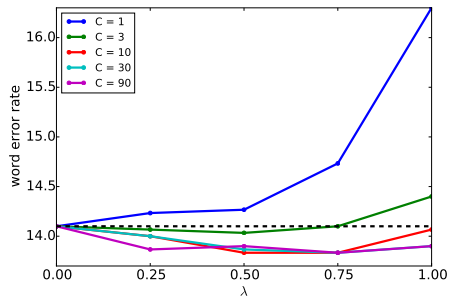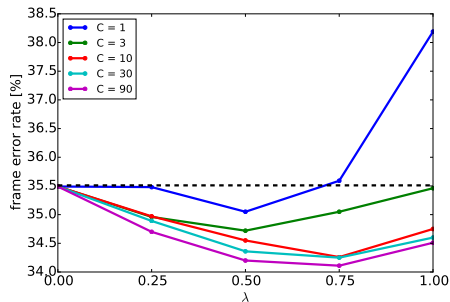


We only keep predictions for classes covering 99% probability mass, we truncate after 90 classes.

# BLENDING LSTMS INTO CNNS

$$L(\lambda) = \lambda \left[ -\sum_j \sum_c p(c|\mathbf{x}_j) \log q(c|\mathbf{x}_j) \right] + (1-\lambda) \left[ -\sum_j \log q(y_j|\mathbf{x}_j) \right]$$

# BLENDING LSTMS INTO CNNS

$$L(\lambda) = \lambda \left[ -\sum_j \sum_c p(c|\mathbf{x}_j) \log q(c|\mathbf{x}_j) \right] + (1-\lambda) \left[ -\sum_j \log q(y_j|\mathbf{x}_j) \right]$$

# RESULTS

| | **FER** | **WER** | model size | execution time |
|---|---|---|---|---|
| Sainath et al.-style CNN | 37.93% | 15.5 | $\approx$ 75M | $\times$ 0.75 |
| vision-style CNN | 35.51% | 14.1 | $\approx$ 75M | $\times$ 1.0 |
| LSTM | 34.15% | 14.4 | $\approx$ 65M | $\times$ 5.8 |
| LSTM $\rightarrow$ CNN blending | **34.11%** | **13.83** | $\approx$ 75M | $\times$ 1.0 |

# SUMMARY OF MODEL BLENDING

- Speech recognition can be improved a lot by vision-style CNNs.

## SUMMARY OF MODEL BLENDING

- Speech recognition can be improved a lot by vision-style CNNs.
- RNNs and CNNs learn different aspects of the data.

- Speech recognition can be improved a lot by vision-style CNNs.
- RNNs and CNNs learn different aspects of the data.
- Recurrent networks for speech recognition may not need to be recurrent.

## SUMMARY OF MODEL BLENDING

- ▶ Speech recognition can be improved a lot by vision-style CNNs.
- ▶ RNNs and CNNs learn different aspects of the data.
- ▶ Recurrent networks for speech recognition may not need to be recurrent.
- ▶ Only "dim knowledge" necessary.

# SELF-COMPRESSION

- Train model A.

## SELF-COMPRESSION

- Train model A.
- Train an indentical model B, mimicking model A with $\lambda = 0.5$.

# SELF-COMPRESSION

- Train model A.
- Train an indentical model B, mimicking model A with $\lambda = 0.5$.
- Model B is more accurate than model A in FER!
    - FER: $35.51 \rightarrow 34.61$.
    - WER: $14.1 \rightarrow 14.1$.

# Thank you!

Do Deep Convolutional Nets Really Need to be Deep (Or Even Convolutional)? ICLR 2017

Gregor Urban, Krzysztof J. Geras, Samira Ebrahimi Kahou, Ozlem Aslan, Shengjie Wang, Rich Caruana,

Abdel-rahman Mohamed, Matthai Philipose, Matthew Richardson

Blending LSTMs into CNNs. ICLR 2016

Krzysztof J. Geras, Abdel-rahman Mohamed, Rich Caruana, Gregor Urban, Shengjie Wang, Ozlem Aslan,

Matthai Philipose, Matthew Richardson and Charles Sutton