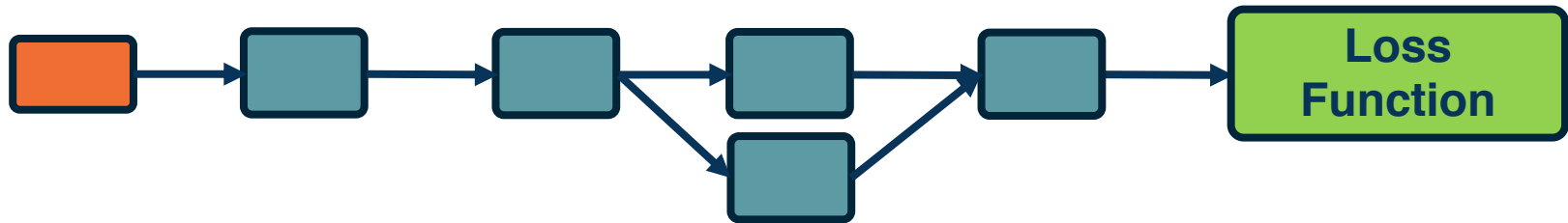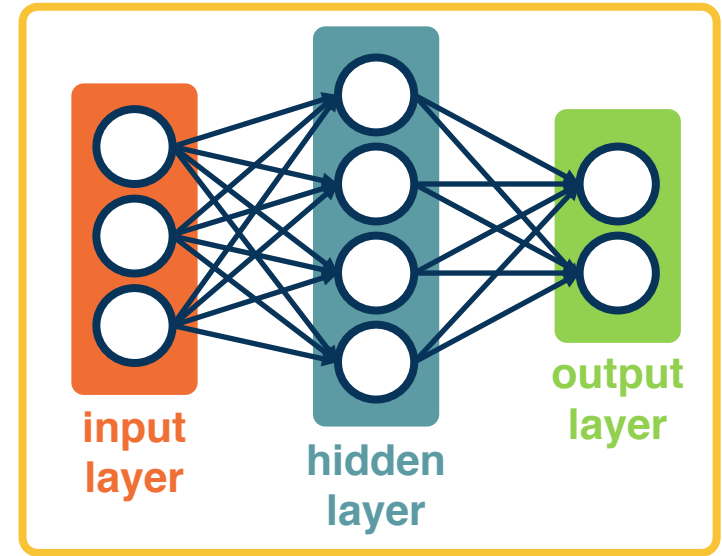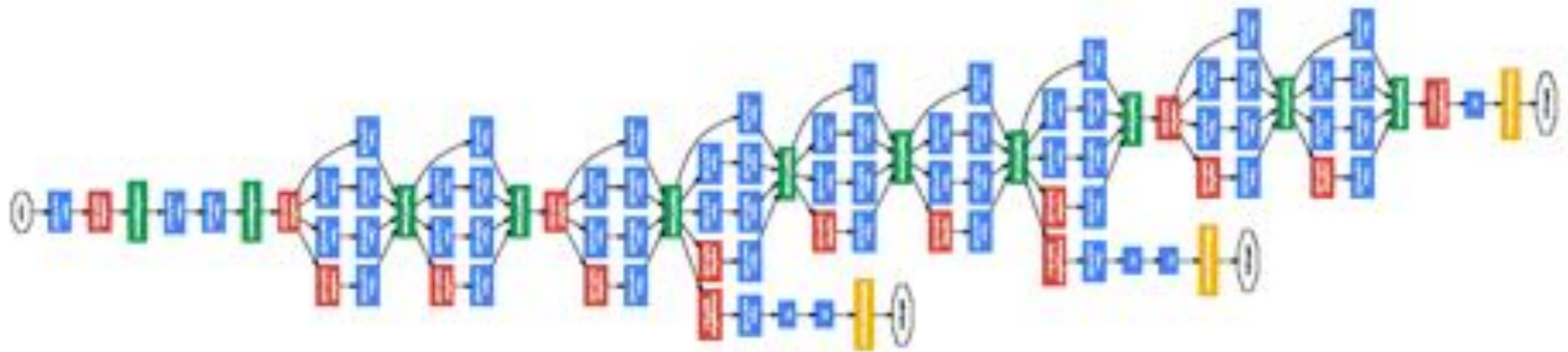**Structures and Structured Representations**

We have seen how **to build and optimize deep feedforward architectures** consisting of linear & non-linear (e.g. ReLU) layers

⬗ This can be generalized to **arbitrary computation graphs**

⬗ **Backpropagation and automatic differentiation** can be used to optimize all parameters via **gradient descent**

FC

Conv
1x1+1(S)

MaxPool
3x3+1(S)

SoftmaxActivation

From: Szegedy et al. Going deeper with convolutions
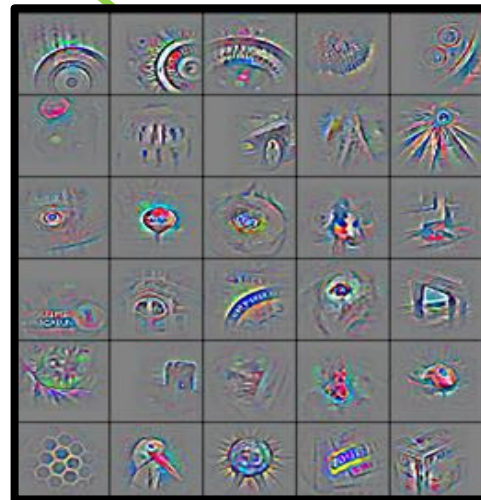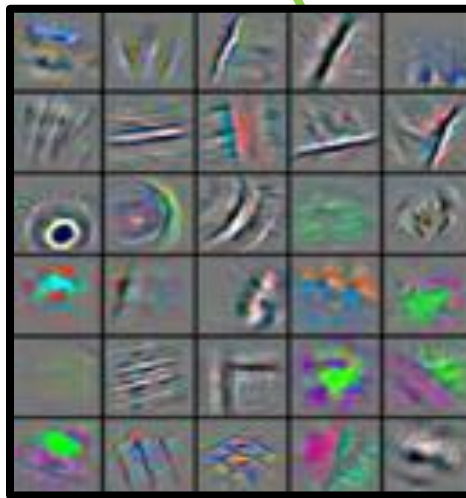
**Convolutional Neural Networks**

Georgia
Tech

Low-Level Feature → Mid-Level Feature → High-Level Feature → Trainable Classifier → "car"

Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

Adapted from slides by: Marc'Aurelio Ranzato, Yann LeCun

**Deep Learning = Hierarchical Compositionality**

Georgia Tech

**VISION**

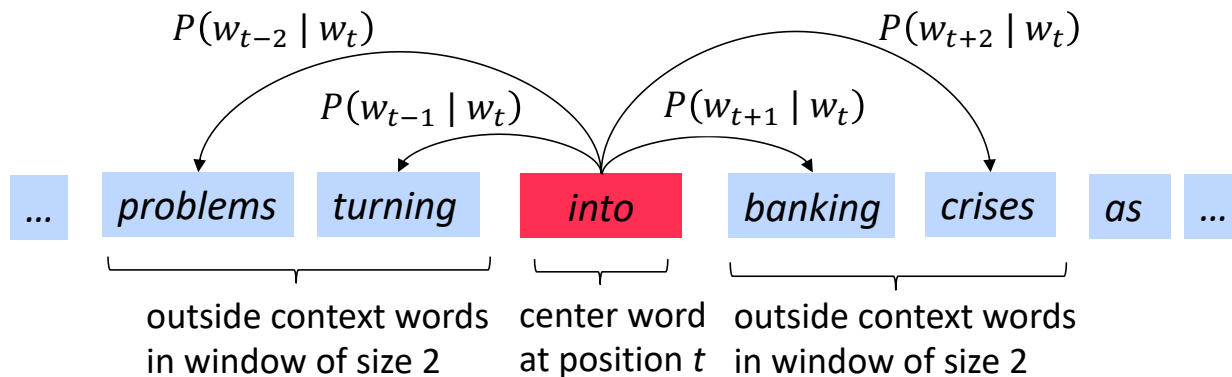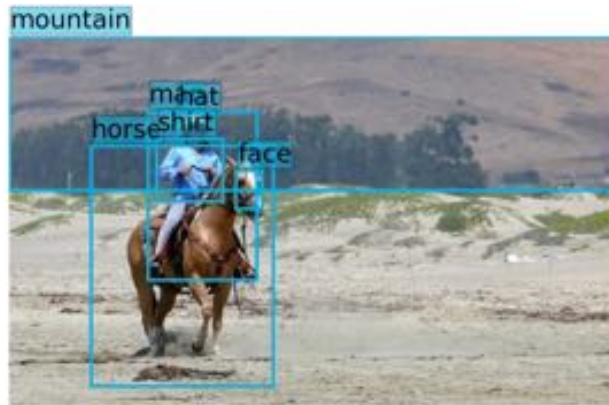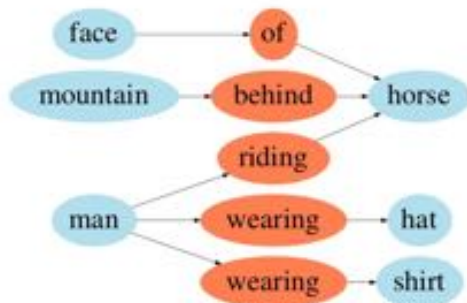pixels ➡ edge ➡ texton ➡ motif ➡ part ➡ object

**SPEECH**

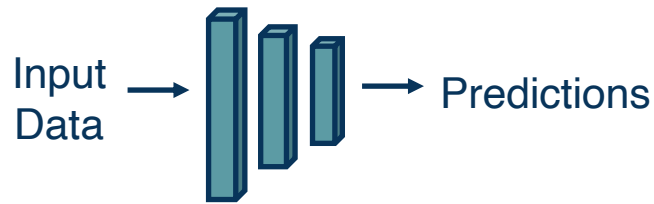sample ➡ spectral band ➡ formant ➡ motif ➡ phone ➡ word

**NLP**

character ➡ word ➡ NP/VP/.. ➡ clause ➡ sentence ➡ story

*Adapted from slides by: Marc'Aurelio Ranzato, Yann LeCun*

**Hierarchical Compositionality**

Georgia Tech

$P(w_{t-2} \mid w_t)$

$P(w_{t+2} \mid w_t)$

$P(w_{t-1} \mid w_t)$

$P(w_{t+1} \mid w_t)$

| ... | *problems* | *turning* | *into* | *banking* | *crises* | *as* | ... |

outside context words in window of size 2
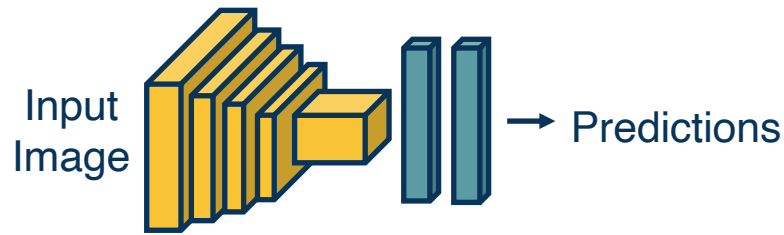
center word at position $t$

outside context words in window of size 2

*Xu et al., "Scene Graph Generation by Iterative Message Passing", 2017*

**Relationships are Important**

Georgia Tech

**Embedding Wikidata Graph [Lerer et al. 19']**

*https://github.com/facebookresearch/ PyTorch-BigGraph*

**Relationships are Everywhere**

Georgia Tech

Input Data → Predictions

**Fully Connected Neural Networks**

Input Image → Predictions
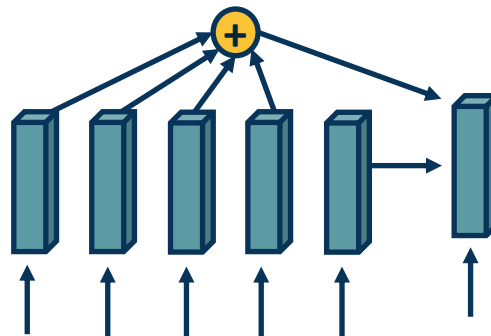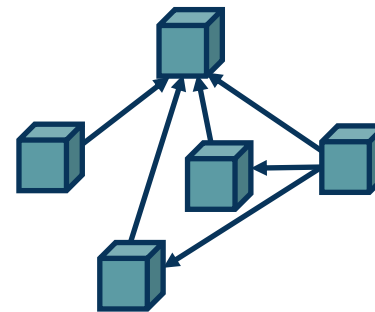
**Convolutional Neural Networks**

**Recurrent Neural Networks**

**Attention-Based Networks**

**Graph-Based Networks**

**The Space of Architectures**

Georgia Tech

**A Multi-Relation Graph**

**Embedding:** A learned map from entities to vectors of numbers that encodes similarity
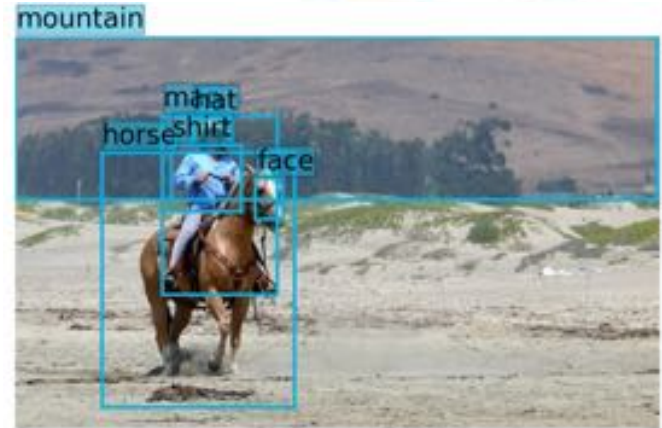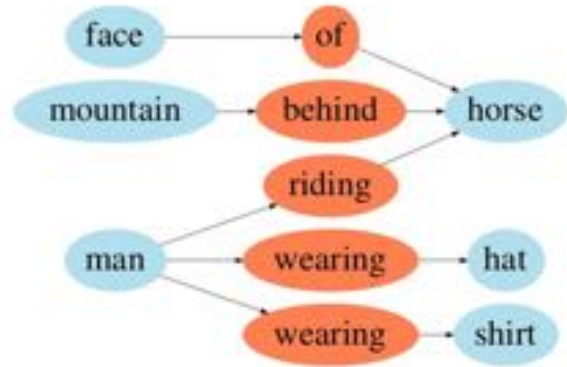
- ⬡ Word embeddings: word ➡ vector

- ⬡ Graph embeddings: node ➡ vector

**Graph Embedding:** Optimize the objective that **connected nodes have more similar embeddings** than unconnected nodes via gradient descent.

*Slide Credit: Adam Lerer*

**Graph Embeddings**

Georgia Tech

When representing structured information, several things are important:

⬡ **State:** Compactly representing all the data we have processed thus far

⬡ **"Neighborhoods"**: What other elements to incorporate?

⬡ Can be seen as selecting from a set of elements

⬡ Typically done with some similarity measure or attention

⬡ **Propagation of information**: How to update information given selected elements



*Slide Credit: Adam Lerer*

**Propagating Information**

- Given a set of vectors $\{u_1, ..., u_N\}$ and a "query" vector $q$
- We can select the most similar vector to $q$ via $p = Softmax(Uq)$



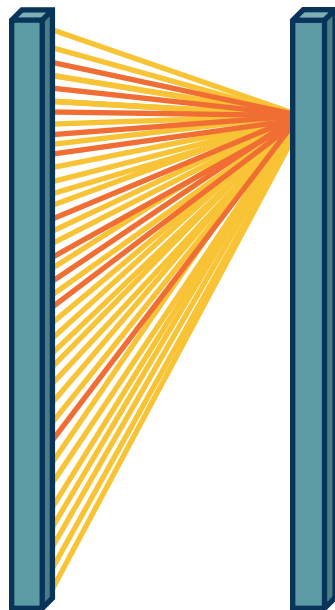$$a_j = \frac{e^{u_j \cdot q}}{\sum_k e^{u_k \cdot q}} \qquad \text{output} = \sum_k a_k u_k$$

Georgia Tech

$$y_i = \frac{1}{C(x)} \sum_{\forall j} f(x_i, x_j) g(x_j)$$

- Output pixel $i$, and $j$ representing all possible input positions

- $f$ is similarity function, $g$ is computes representation of input element $j$

- **Examples:**

$$f(x_i, x_j) = e^{x_i^T x_j}$$

$$g(x_j) = W_g x_j$$

*Wang et al., "Non-local Neural Networks", 2017*

1024 x 1024 Pixel Image

~1M element Vector **(M)**

Non-Local Layer

**Example: Non-Local Neural Networks**

Georgia Tech

$$y_i = \frac{1}{C(x)} \sum_{\forall j} f(x_i, x_j) g(x_j)$$

◆ Output pixel $i$, and $j$ representing all possible input positions

◆ $f$ is similarity function, $g$ is computes representation of input element $j$

◆ **Examples:**
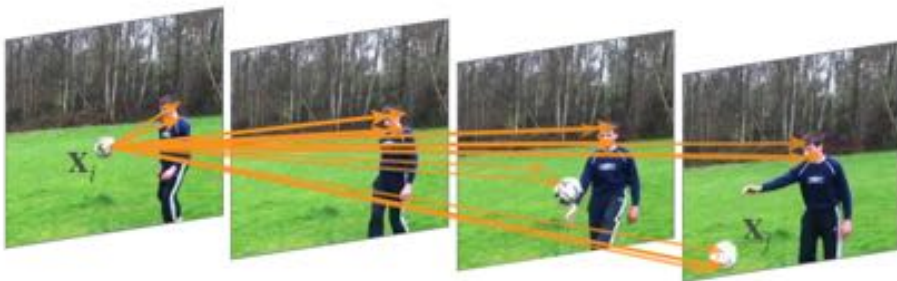
$$f(x_i, x_j) = e^{x_i^T x_j}$$

$$g(x_j) = W_g x_j$$



Figure 1. A spacetime *non-local* operation in our network trained for video classification in Kinetics. A position $\mathbf{x}_i$'s response is computed by the weighted average of the features of *all* positions $\mathbf{x}_j$ (only the highest weighted ones are shown here). In this example computed by our model, note how it relates the ball in the first frame to the ball in the last two frames. More examples are in Figure 3.

*Wang et al., "Non-local Neural Networks", 2017*

**Example: Non-Local Neural Networks**

Georgia Tech