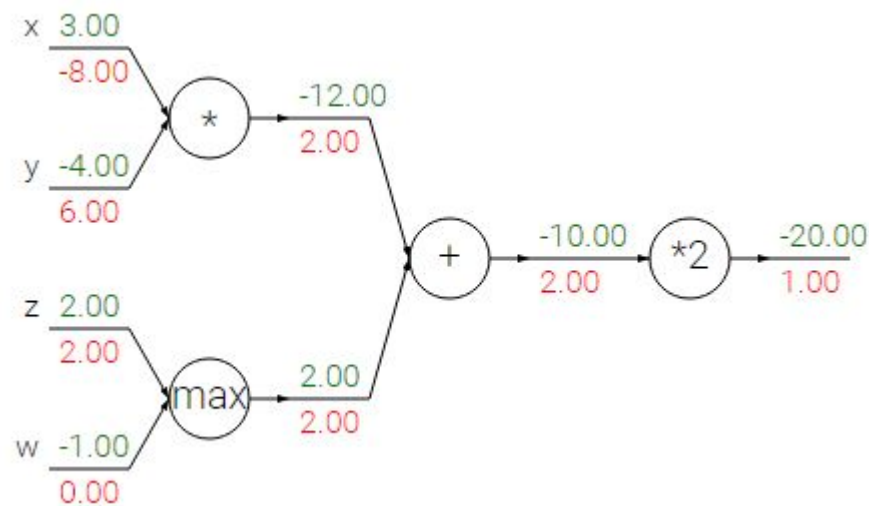# Making Neural Networks Work Well!
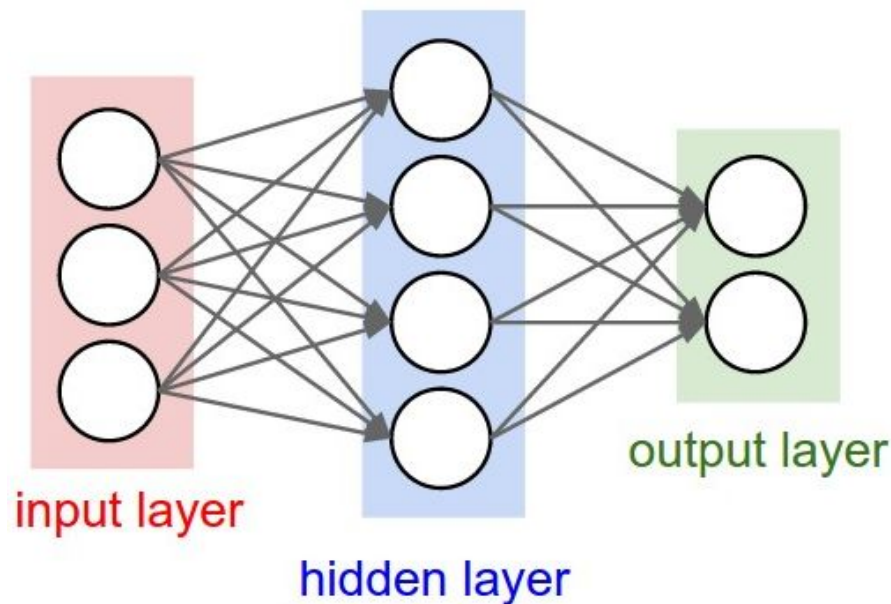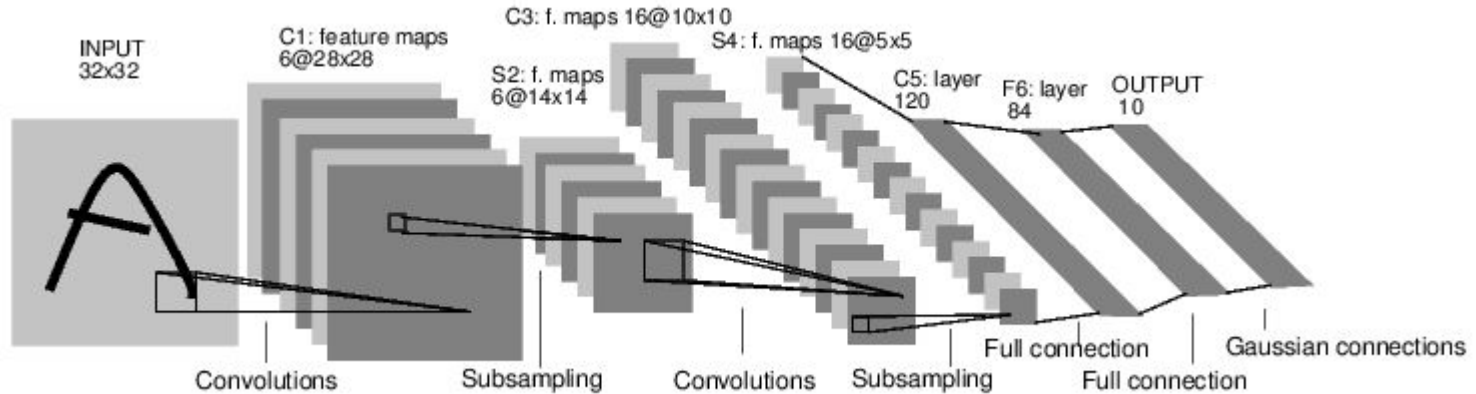
By: Ismail Elezi

ismail.elezi@gmail.com

# Backpropagation without formulas
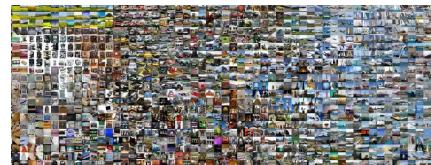


input layer

hidden layer

output layer

# Convolutional Neural Networks (CNN)

# But we all know that (since the end of eighties)

1. We have more data - from Lena to ImageNet.



2. We have more computing power, GPUs are really good at this.



3. Last but not least, we have new ideas (and this is the focus of this lecture).
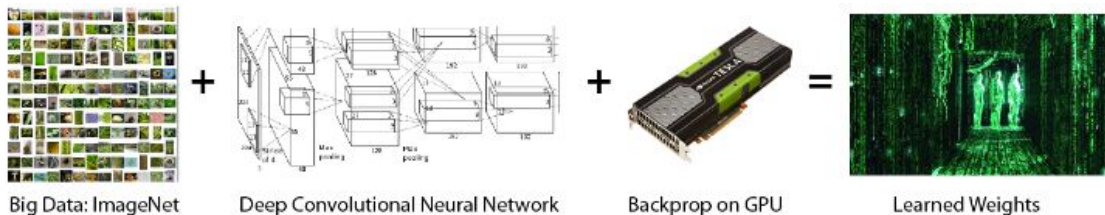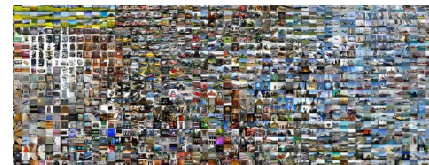
# But we all know that (since the end of eighties)

1. We have more data - from Lena to ImageNet.



2. We have more computing power, GPUs are really good at this.



3. Last but not least, we have new ideas (and this is the focus of this lecture).





Big Data: ImageNet + Deep Convolutional Neural Network + Backprop on GPU = Learned Weights

# 1) Use Mini-Batches

# 2) Softmax and Cross Entropy



$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

unnormalized probabilities

| | | | |
|---|---|---|---|
| cat | **3.2** | **24.5** | **0.13** |
| car | 5.1 | 164.0 | 0.87 |
| frog | -1.7 | 0.18 | 0.00 |

exp → normalize →

L_i = -log(0.13) = 0.89

unnormalized log probabilities

probabilities

# 3) Use Rectified Linear Units (ReLU)



$$S(t) = \frac{1}{1 + e^{-t}}.$$

$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$f(x) = \max(0, x)$$

# 4) Consider using dropout (and L2 regularization)

- A new form of regularization ([Srivastava et al., 2014](#)).


- Could be interpreted as an ensemble neural network.



- That doesn't mean that you shouldn't also use L2 regularization.



(a) Standard Neural Net

(b) After applying dropout.

# 5) Use Batch-Normalization

$$\textbf{Input:} \quad \text{Values of } x \text{ over a mini-batch: } \mathcal{B} = \{x_{1...m}\};$$
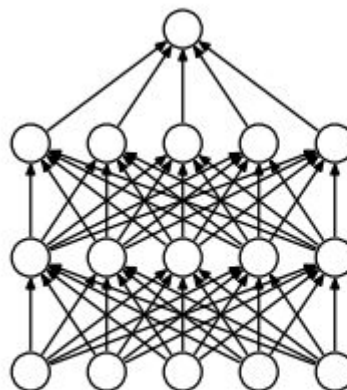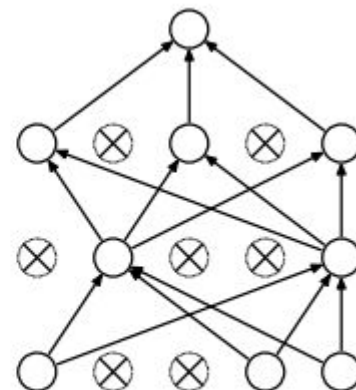$$\qquad\qquad \text{Parameters to be learned: } \gamma, \beta$$
$$\textbf{Output:} \quad \{y_i = \text{BN}_{\gamma,\beta}(x_i)\}$$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m}\sum_{i=1}^{m} x_i \qquad\qquad \text{// mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m}\sum_{i=1}^{m}(x_i - \mu_{\mathcal{B}})^2 \qquad \text{// mini-batch variance}$$

$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad\qquad \text{// normalize}$$

$$y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i) \qquad \text{// scale and shift}$$

- Do normalization using mean and standard deviation on each channel of the image.

- Do scaling and shifting.

- Find the parameters of scaling and shifting using backpropagation (this is really awesome, why other people didn't think before about it?!).

- Read the paper (Ioffe and Szegedy) :)

# 6) Different Optimizers - use Adam :)

- Stochastic gradient descent is explained in every machine learning book (and course).
- But, no one uses it on vanilla format.
- And adding just momentum is so last century.
- Nesterov's momentum, Adagrad, Adadelta and RMSProp are all very good choices.
- Adam (Kingma and Ba) is an excellent choice.
- Second order optimization algorithms don't work well.
  $$x \leftarrow x - [Hf(x)]^{\wedge}(-1) \nabla f(x)$$

# 7) Searching for the best hyperparameters

- The number of hyperparameters is very large on neural networks ( # hyperparameters >> 10).
- There are a lot of rules of thumb for some of the hyperparameters, and intuition for many others.
- Grid search is totally unfeasible, random search (Bergstra and Bengio, 2012) works much better.

# 8) Testing set is not enough anymore (if it ever was)!

1. In some algorithms you might be able to use just a training set and a testing set.

2. You might think to do that in ANN.

3. And then you will overfit your testing set, and so your reported results will be fake.

4. Use a validation set.

5. Pedro Domingos' paper is a must read on this (and many other useful things to know, extra points for being easy to read)



Training
Validation
Testing

# 9) Use transfer learning - whenever is possible

- Modern CNNs take several weeks to be trained in a cluster of GPUs.
- Instead of training a CNN from the scratch, download a pretrained one, remove the last layer, and finetune it.
- This not only it saves a lot of precious time, but it can also give you better results.
- Morale of the story: Use transfer learning in CNNs.

# 10) Don't be scared from large models



Credit: C. Zhang

# Lots and lots more

Andrej Karpathy's notes are the best tutorial you can hope to find on neural networks (and the videos of the course are available).

And now we have a book in deep learning (Goodfellow, Bengio and Courville, Chapters 8 and 11).

Andrew Ng's book (still on draft stage) is a nice book on how to use machine learning and deep learning efficiently. His book is very much based on his presentations.

Don't forget to implement all the algorithms I mentioned in this presentation. The best way of learning algorithms is by implementing them.

# Recurrent Neural Networks

By: Ismail Elezi

ismail.elezi@gmail.com

# RNN intro

Recurrent Neural Networks are Turing complete, yipes!

Remember from the class the universal approximation theorem?!

This combination looks kind of powerful, right?

**If training vanilla neural nets is optimization over functions, training recurrent nets is optimization over programs.**

# RNN intro



one to one | one to many | many to one | many to many | many to many

Credit: Andrej Karpathy

# Vanilla RNN



Recurrent Neural Networks have loops.

An unrolled recurrent neural network.

Credit: Christopher Olah

# Vanilla RNN (2)



$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t)$$

$$y_t = W_{hy} h_t$$

```python
class RNN:
    # ...
    def step(self, x):
        # update the hidden state
        self.h = np.tanh(np.dot(self.W_hh, self.h) + np.dot(self.W_xh, x))
        # compute the output vector
        y = np.dot(self.W_hy, self.h)
        return y
rnn = RNN()
y = rnn.step(x) # x is an input vector, y is the RNN's output vector
```

Credit: Goodfellow et al.

# The vanishing gradient problem

$$h_t = \theta\phi(h_{t-1}) + \theta_x x_t$$

$$y_t = \theta_y\phi(h_t)$$

$$\frac{\partial E}{\partial \theta} = \sum_{t=1}^{s} \frac{\partial E_t}{\partial \theta}$$

$$\frac{\partial E}{\partial \theta} = \sum_{t=1}^{s} \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial \theta}$$

# The vanishing gradient problem (2)

$$\frac{\partial E}{\partial \theta} = \sum_{t=1}^{s} \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial \theta}$$

$$\frac{\partial h_t}{\partial h_k} = \prod_{i=k+1}^{t} \frac{\partial h_i}{\partial h_{i-1}} = \prod_{i=k+1}^{t} \theta^T diag[\phi^{'}(h_{i-1})]$$

$$\left\| \frac{\partial h_i}{\partial h_{i-1}} \right\| \leq \left\| \theta^T \right\| \left\| diag[\phi^{'}(h_{i-1})] \right\| \leq \gamma_\theta \gamma_\phi$$

$$\left\| \frac{\partial h_t}{\partial h_k} \right\| \leq (\gamma_\theta \gamma_\phi)^{t-k}$$

# Long Short Term Memory (LSTM) - looks scary

# But they really are not



The repeating module in a standard RNN contains a single layer.

The repeating module in an LSTM contains four interacting layers.

Credit: Christopher Olah

# LSTM (1)



$$f_t = \sigma\left(W_f \cdot [h_{t-1}, x_t] + b_f\right)$$

Credit: Christopher Olah

# LSTM (2)



$$i_t = \sigma\left(W_i \cdot [h_{t-1}, x_t] + b_i\right)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Credit: Christopher Olah

# LSTM (3)



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Credit: Christopher Olah

# LSTM (4)



$$o_t = \sigma \left( W_o \left[ h_{t-1}, x_t \right] + b_o \right)$$

$$h_t = o_t * \tanh \left( C_t \right)$$

Credit: Christopher Olah

# LSTM in code

```python
def lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b):
    """
    Forward pass for a single timestep of an LSTM.

    The input data has dimension D, the hidden state has dimension H, and we use
    a minibatch size of N.

    Inputs:
    - x: Input data, of shape (N, D)
    - prev_h: Previous hidden state, of shape (N, H)
    - prev_c: previous cell state, of shape (N, H)
    - Wx: Input-to-hidden weights, of shape (D, 4H)
    - Wh: Hidden-to-hidden weights, of shape (H, 4H)
    - b: Biases, of shape (4H,)

    Returns a tuple of:
    - next_h: Next hidden state, of shape (N, H)
    - next_c: Next cell state, of shape (N, H)
    - cache: Tuple of values needed for backward pass.
    """
    next_h, next_c, cache = None, None, None

    N, H = prev_h.shape
    # 1
    a = np.dot(x, Wx) + np.dot(prev_h, Wh) + b

    # 2
    ai = a[:, :H]
    af = a[:, H:2*H]
    ao = a[:, 2*H:3*H]
    ag = a[:, 3*H:]

    # 3
    i = sigmoid(ai)
    f = sigmoid(af)
    o = sigmoid(ao)
    g = np.tanh(ag)

    # 4
    next_c = f * prev_c + i * g

    # 5
    next_h = o * np.tanh(next_c)

    cache = i, f, o, g, a, ai, af, ao, ag, Wx, Wh, b, prev_h, prev_c, x, next_c, next_h

    return next_h, next_c, cache
```

```python
def lstm_step_backward(dnext_h, dnext_c, cache):
    """
    Backward pass for a single timestep of an LSTM.

    Inputs:
    - dnext_h: Gradients of next hidden state, of shape (N, H)
    - dnext_c: Gradients of next cell state, of shape (N, H)
    - cache: Values from the forward pass

    Returns a tuple of:
    - dx: Gradient of input data, of shape (N, D)
    - dprev_h: Gradient of previous hidden state, of shape (N, H)
    - dprev_c: Gradient of previous cell state, of shape (N, H)
    - dWx: Gradient of input-to-hidden weights, of shape (D, 4H)
    - dWh: Gradient of hidden-to-hidden weights, of shape (H, 4H)
    - db: Gradient of biases, of shape (4H,)
    """
    dx, dh, dc, dWx, dWh, db = None, None, None, None, None, None

    i, f, o, g, a, ai, af, ao, ag, Wx, Wh, b, prev_h, prev_c, x, next_c, next_h = cache

    # backprop into step 5
    do = np.tanh(next_c) * dnext_h
    dnext_c += o * (1 - np.tanh(next_c) ** 2) * dnext_h

    # backprop into 4
    df = prev_c * dnext_c
    dprev_c = f * dnext_c
    di = g * dnext_c
    dg = i * dnext_c

    # backprop into 3
    dai = sigmoid(ai) * (1 - sigmoid(ai)) * di
    daf = sigmoid(af) * (1 - sigmoid(af)) * df
    dao = sigmoid(ao) * (1 - sigmoid(ao)) * do
    dag = (1 - np.tanh(ag) ** 2) * dg

    # backprop into 2
    da = np.hstack((dai, daf, dao, dag))

    # backprop into 1
    db = np.sum(da, axis = 0)
    dprev_h = np.dot(Wh, da.T).T
    dWh = np.dot(prev_h.T, da)
    dx = np.dot(da, Wx.T)
    dWx = np.dot(x.T, da)

    return dx, dprev_h, dprev_c, dWx, dWh, db
```
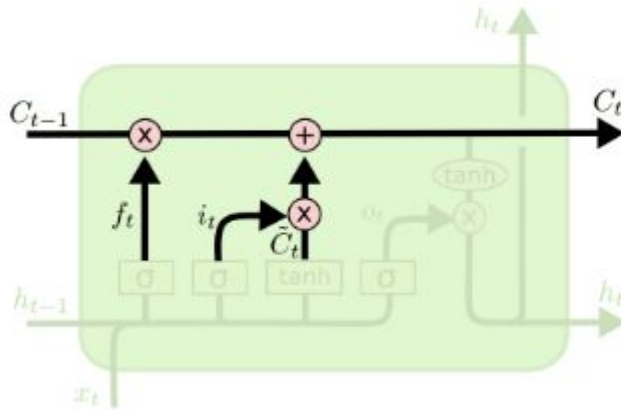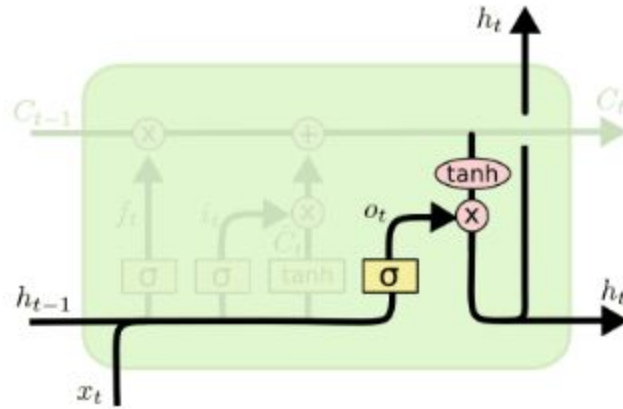
# RNN handwriting generation demo

# A Song of Long and Short Term Memory

Thou will fly not her more beyond his feet.
A look at him, worth it went sunlight. The handbell was no days at lasquid.
The fellow. "The knees. The hair was vest low and touch closer to reterming the vowmen."
And the notion, and one is in the bloodriders and horring to cheer of praying a teur women, I land and looking stares, certain. I bear Bran told me liked that hawers came the
queen could be help you on it. He won their bannervles than hard two. **Jon had died**.

Undered wagons before he could hear the hills. Cuse behind them down from Harrenhal, had given day atop. While, but everya, and the lake. Stannis, her shoulder. Have them spinning with properior of raisings or a stable, Mooda'an of Lord Tywin, anyone said she could then he took twelve from some helm of the king's father would have
seen themanly helpless and hatish with the rear; an age to go up, feed him and limp and clawing through sheep of fight.

It's any of his head againand, but
Lord Tywin raping for Westeros. The king's chance, Iron. "Where he wanted so bride his own words from off insists and brail in empty food," her crops: dreaple white and woman
spounded onto her place. "Little fear?" The course pointed over a rock and life. They were done, Lady Stark like whatsa think to die seen

Could-becaping a bride of her wolfstorm, you're boils as he passed down.
Bronn still guess and gave her every thinking. "I have only warm everything, it sellsboads," he wondered what she turned his clung butter ice husband to emerally now came came line to Needle. Jober. She is true, and the High Mallyrion of the again, we need to instear up the
sellsword, as contemplate back.
One of their horses: Seep, north his beard and when so langess, rangers spray in the rose. There around him.

Blhood Scollop of his son think the tail forgiver and don't be sure, and kinds your
Wees. I was, walls, I know your sisters. Maester Luwin battle!"
"I'd ice. The lovely in thick?" she asked Jhiquire seemed; to the needles remainne in the risk levies
now, my lord like wall and left it. After sail as a servA MUCKENY, a guards crawl stoncelfing north to hear more so one belts so mind out of the firell wondered someone said that wench to surver cut down the notion.

# War and Peace

tyntd-iafhatawiaoihrdemot  lytdws  e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e
plia tklrgd t o idoe ns,smtt   h ne etie h,hregtrs nigtike,aoaenns lng

"Tmont thithey" fomesscerliund
Keushey. Thom here
sheulke, anmerenith ol sivh I lalterthend Bleipile shuwy fil on aseterlome
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."

we counter. He stutn co des. His stanted out one ofler that concossions and was
to gearang reay Jotrets and with fre colt otf paitt thin wall. Which das stimn

Aftair fall unsuch that the hall for Prince Velzonski's that me of
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort
how, and Gogition is so overelical and ofter.

"Kite vouch!" he repeated by her
door. "But I would be done and quarts, feeling, then, son is people...."

"Why do what that day," replied Natasha, and wishing to himself the fact the
princess, Princess Mary was easier, fed in had oftened him.
Pierre aking his soul came to the packs and drove up his father-in-law women.

Credit: A. Karpathy

# Faking mathematics

For $\bigoplus_{n=1,\ldots,m}$ where $\mathcal{L}_{m_\bullet} = 0$, hence we can find a closed subset $\mathcal{H}$ in $\mathcal{H}$ and any sets $\mathcal{F}$ on $X$, $U$ is a closed immersion of $S$, then $U \to T$ is a separated algebraic space.

*Proof.* Proof of (1). It also start we get

$$S = \mathrm{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by $\coprod Z \times_U U \to V$. Consider the maps $M$ along the set of points $Sch_{fppf}$ and $U \to U$ is the fibre category of $S$ in $U$ in Section, ?? and the fact that any $U$ affine, see Morphisms, Lemma ??. Hence we obtain a scheme $S$ and any open subset $W \subset U$ in $Sh(G)$ such that $\mathrm{Spec}(R') \to S$ is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that $f_i$ is of finite presentation over $S$. We claim that $\mathcal{O}_{X,x}$ is a scheme where $x, x', s'' \in S'$ such that $\mathcal{O}_{X,x'} \to \mathcal{O}'_{X',x'}$ is separated. By Algebra, Lemma ?? we can define a map of complexes $\mathrm{GL}_{S'}(x'/S'')$ and we win. $\square$

To prove study we see that $\mathcal{F}|_U$ is a covering of $\mathcal{X}'$, and $\mathcal{T}_i$ is an object of $\mathcal{F}_{X/S}$ for $i > 0$ and $\mathcal{F}_p$ exists and let $\mathcal{F}_i$ be a presheaf of $\mathcal{O}_X$-modules on $\mathcal{C}$ as a $\mathcal{F}$-module. In particular $\mathcal{F} = U/\mathcal{F}$ we have to show that

$$\widetilde{M}^\bullet = \mathcal{I}^\bullet \otimes_{\mathrm{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1}\mathcal{F})$$

is a unique morphism of algebraic stacks. Note that

$$\mathrm{Arrows} = (Sch/S)_{fppf}^{opp}, (Sch/S)_{fppf}$$

and

$$V = \Gamma(S, \mathcal{O}) \longmapsto (U, \mathrm{Spec}(A))$$

is an open subset of $X$. Thus $U$ is affine. This is a continuous map of $X$ is the inverse, the groupoid scheme $S$.

*Proof.* See discussion of sheaves of sets. $\square$

The result for prove any open covering follows from the less of Example ??. It may replace $S$ by $X_{spaces,\acute{e}tale}$ which gives an open subspace of $X$ and $T$ equal to $S_{Zar}$, see Descent, Lemma ??. Namely, by Lemma ?? we see that $R$ is geometrically regular over $S$.

**Lemma 0.1.** *Assume (3) and (3) by the construction in the description.*

*Suppose $X = \lim |X|$ (by the formal open covering $X$ and a single map $\underline{Proj}_X(\mathcal{A}) = \mathrm{Spec}(B)$ over $U$ compatible with the complex*

$$Set(\mathcal{A}) = \Gamma(X, \mathcal{O}_{X,\mathcal{O}_X}).$$

*When in this case of to show that $\mathcal{Q} \to \mathcal{C}_{Z/X}$ is stable under the following result in the second conditions of (1), and (3). This finishes the proof. By Definition ?? (without element is when the closed subschemes are catenary. If $T$ is surjective we may assume that $T$ is connected with residue fields of $S$. Moreover there exists a closed subspace $Z \subset X$ of $X$ where $U$ in $X'$ is proper (some defining as a closed subset of the uniqueness it suffices to check the fact that the following theorem*

(1) $f$ *is locally of finite type. Since $S = \mathrm{Spec}(R)$ and $Y = \mathrm{Spec}(R)$.*

*Proof.* This is form all sheaves of sheaves on $X$. But given a scheme $U$ and a surjective étale morphism $U \to X$. Let $U \cap U = \coprod_{i=1,\ldots,n} U_i$ be the scheme $X$ over $S$ at the schemes $X_i \to X$ and $U = \lim_i X_i$. $\square$

The following lemma surjective restrocomposes of this implies that $\mathcal{F}_{x_0} = \mathcal{F}_{x_0} = \mathcal{F}_{X,\ldots,0}$.

**Lemma 0.2.** *Let $X$ be a locally Noetherian scheme over $S$, $E = \mathcal{F}_{X/S}$. Set $\mathcal{I} = \mathcal{J}_1 \subset \mathcal{I}'_n$. Since $\mathcal{I}^n \subset \mathcal{I}^n$ are nonzero over $i_0 \leq \mathfrak{p}$ is a subset of $\mathcal{J}_{n,0} \circ \overline{A}_2$ works.*

**Lemma 0.3.** *In Situation ??. Hence we may assume $\mathfrak{q}' = 0$.*

*Proof.* We will use the property we see that $\mathfrak{p}$ is the mext functor (??). On the other hand, by Lemma ?? we see that

$$D(\mathcal{O}_{X'}) = \mathcal{O}_X(D)$$

where $K$ is an $F$-algebra where $\delta_{n+1}$ is a scheme over $S$. $\square$

Credit: A. Karpathy

# Faking mathematics (2)

*Proof.* Omitted. □

**Lemma 0.1.** *Let* $\mathcal{C}$ *be a set of the construction.*

*Let* $\mathcal{C}$ *be a gerber covering. Let* $\mathcal{F}$ *be a quasi-coherent sheaves of* $\mathcal{O}$-*modules. We have to show that*

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

.

*Proof.* This is an algebraic space with the composition of sheaves $\mathcal{F}$ on $X_{\acute{e}tale}$ we have

$$\mathcal{O}_X(\mathcal{F}) = \{morph_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where $\mathcal{G}$ defines an isomorphism $\mathcal{F} \to \mathcal{F}$ of $\mathcal{O}$-modules. □

**Lemma 0.2.** *This is an integer* $\mathcal{Z}$ *is injective.*

*Proof.* See Spaces, Lemma ??. □

**Lemma 0.3.** *Let* $S$ *be a scheme. Let* $X$ *be a scheme and* $X$ *is an affine open covering. Let* $\mathcal{U} \subset \mathcal{X}$ *be a canonical and locally of finite type. Let* $X$ *be a scheme. Let* $X$ *be a scheme which is equal to the formal complex.*

*The following to the construction of the lemma follows.*

*Let* $X$ *be a scheme. Let* $X$ *be a scheme covering. Let*

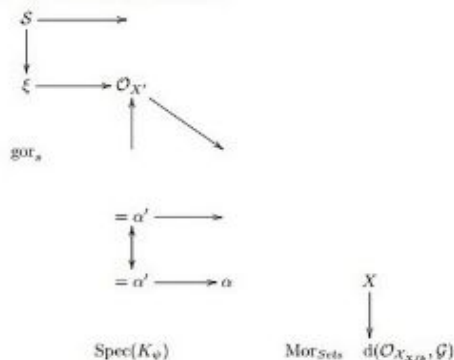$$b: X \to Y' \to Y \to Y \to Y' \times_X Y \to X.$$

*be a morphism of algebraic spaces over* $S$ *and* $Y$.

*Proof.* Let $X$ be a nonzero scheme of $X$. Let $X$ be an algebraic space. Let $\mathcal{F}$ be a quasi-coherent sheaf of $\mathcal{O}_X$-modules. The following are equivalent

(1) $\mathcal{F}$ is an algebraic space over $S$.
(2) If $X$ is an affine open covering.

Consider a common structure on $X$ and $X$ the functor $\mathcal{O}_X(U)$ which is locally of finite type. □

This since $\mathcal{F} \in \mathcal{F}$ and $x \in \mathcal{G}$ the diagram



is a limit. Then $\mathcal{G}$ is a finite type and assume $S$ is a flat and $\mathcal{F}$ and $\mathcal{G}$ is a finite type $f_*$. This is of finite type diagrams, and

- the composition of $\mathcal{G}$ is a regular sequence.
- $\mathcal{O}_{X'}$ is a sheaf of rings. □

*Proof.* We have see that $X = \mathrm{Spec}(R)$ and $\mathcal{F}$ is a finite type representable by algebraic space. The property $\mathcal{F}$ is a finite morphism of algebraic stacks. Then the cohomology of $X$ is an open neighbourhood of $U$. □

*Proof.* This is clear that $\mathcal{G}$ is a finite presentation, see Lemmas ??. A reduced above we conclude that $U$ is an open covering of $\mathcal{C}$. The functor $\mathcal{F}$ is a "field

$$\mathcal{O}_{X,x} \longrightarrow \mathcal{F}_{\overline{x}} \ -1(\mathcal{O}_{X_{\acute{e}tale}}) \longrightarrow \mathcal{O}_{X_\lambda}^{-1} \mathcal{O}_{X_\lambda}(\mathcal{O}_{X_s}^{\mathcal{G}})$$

is an isomorphism of covering of $\mathcal{O}_{X_s}$. If $\mathcal{F}$ is the unique element of $\mathcal{F}$ such that $X$ is an isomorphism.
The property $\mathcal{F}$ is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme $\mathcal{O}_X$-algebra with $\mathcal{F}$ are opens of finite type over $S$. If $\mathcal{F}$ is a scheme theoretic image points. □

If $\mathcal{F}$ is a finite direct sum $\mathcal{O}_{X_\lambda}$ is a closed immersion, see Lemma ??. This is a sequence of $\mathcal{F}$ is a similar morphism.

Credit: A. Karpathy

# This code looks really great

```
/*
 * Increment the size file of the new incorrect UI_FILTER group information
 * of the size generatively.
 */
static int indicate_policy(void)
{
  int error;
  if (fd == MARN_EPT) {
    /*
     * The kernel blank will coeld it to userspace.
     */
    if (ss->segment < mem_total)
      unblock_graph_and_set_blocked();
    else
      ret = 1;
    goto bail;
  }
  segaddr = in_SB(in.addr);
  selector = seg / 16;
  setup_works = true;
  for (i = 0; i < blocks; i++) {
    seq = buf[i++];
    bpf = bd->bd.next + i * search;
    if (fd) {
      current = blocked;
    }
  }
  rw->name = "Getjbbregs";
  bprm_self_clearl(&iv->version);
  regs->new = blocks[(BPF_STATS << info->historidac)] | PFMR_CLOBATHINC_SECONDS << 12;
  return segtable;
}
```

```
static void action_new_function(struct s_stat_info *wb)
{
  unsigned long flags;
  int lel_idx_bit = e->edd, *sys & ~((unsigned long) *FIRST_COMPAT);
  buf[0] = 0xFFFFFFFF & (bit << 4);
  min(inc, slist->bytes);
  printk(KERN_WARNING "Memory allocated %02x/%02x, "
    "original MLL instead\n"),
    min(min(multi_run - s->len, max) * num_data_in),
    frame_pos, sz + first_seg);
  div_u64_w(val, inb_p);
  spin_unlock(&disk->queue_lock);
  mutex_unlock(&s->sock->mutex);
  mutex_unlock(&func->mutex);
  return disassemble(info->pending_bh);
}

static void num_serial_settings(struct tty_struct *tty)
{
  if (tty == tty)
    disable_single_st_p(dev);
  pci_disable_spool(port);
  return 0;
}

static void do_command(struct seq_file *m, void *v)
{
  int column = 32 << (cmd[2] & 0x80);
  if (state)
    cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
  else
    seq = 1;
  for (i = 0; i < 16; i++) {
    if (k & (1 << 1))
      pipe = (in_use & UMXTHREAD_UNCCA) +
        ((count & 0x00000000ffffff8) & 0x000000f) << 8;
    if (count == 0)
      sub(pid, ppc_md.kexec_handle, 0x20000000);
    pipe_set_bytes(i, 0);
  }
```

Credit: A. Karpathy

# Can a computer understand these pictures?



A yellow bus driving down a road with green trees and green grass in the background.

Living room with white couch and blue carpeting. The room in the apartment gets some afternoon sun.

Credit: Andrew Ng

# Image captioning



a man holding a tennis racket on a court <END>
GT:<START> a guy playing tennis waiting for the serve <END>
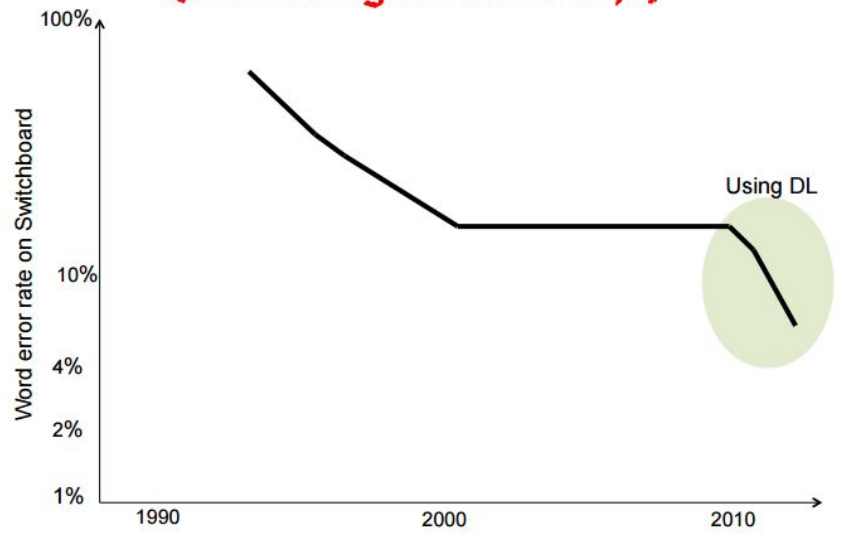


a man is riding a skateboard at a skate park <END>
GT:<START> the young man is carrying his skateboard <END>

The dramatic impact of Deep Learning on Speech Recognition (according to Microsoft)
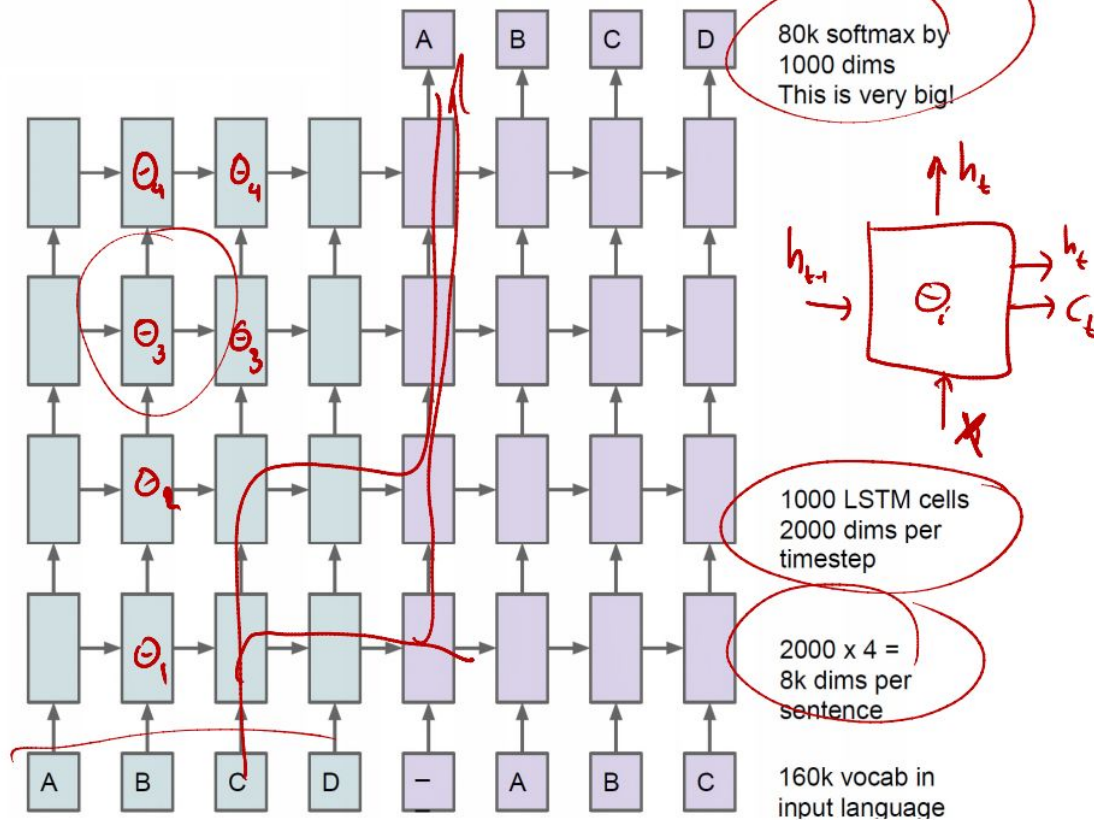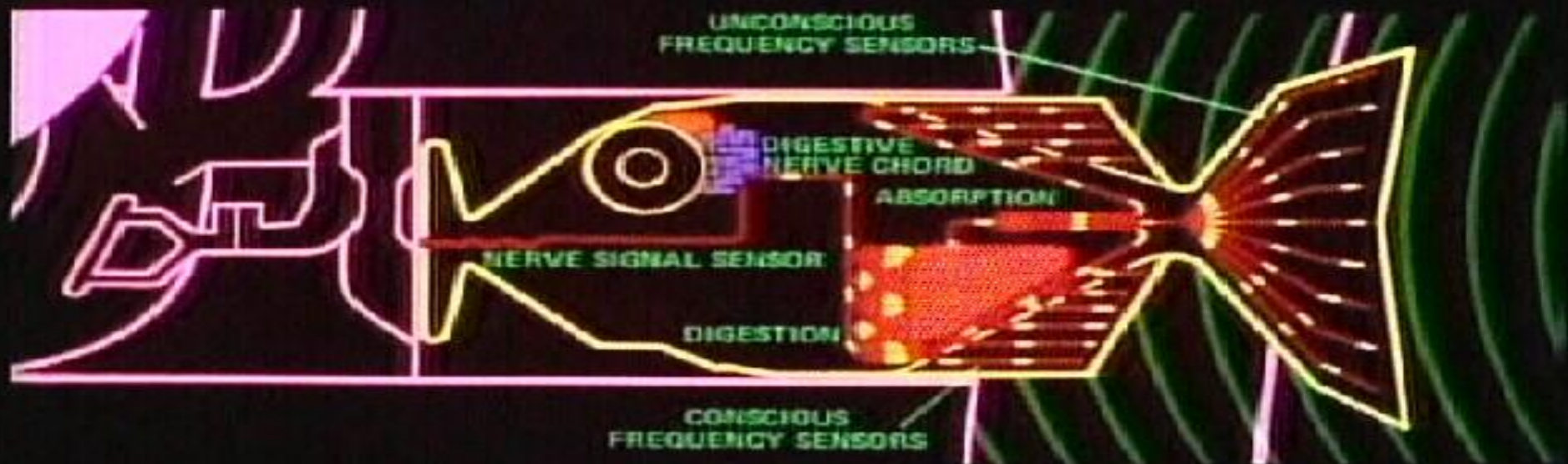
Speech recognition performance

# LSTMs for sequence to sequence prediction



80k softmax by
1000 dims
This is very big!

1000 LSTM cells
2000 dims per
timestep

2000 x 4 =
8k dims per
sentence

160k vocab in
input language

*"Deep Learning waves have lapped at the shores of computational linguistics for several years now, but 2015 seems like the year when the full force of the tsunami hit the major Natural Language Processing (NLP) conferences."*

- *Chris Manning*

Credit: N. de Freitas

# BABEL FISH

UNCONSCIOUS
FREQUENCY SENSORS

DIGESTIVE
NERVE CHORD

ABSORPTION

NERVE SIGNAL SENSOR

DIGESTION

CONSCIOUS
FREQUENCY SENSORS

STICK ONE IN YOUR EAR, YOU CAN INSTANTLY
UNDERSTAND ANYTHING SAID TO YOU IN ANY FORM
OF LANGUAGE: THE SPEECH YOU HEAR DECODES THE
BRAIN WAVE MATRIX.

# CNN and RNN are converging?



Credit: J. Schmidhuber

# Neural Turing Machines
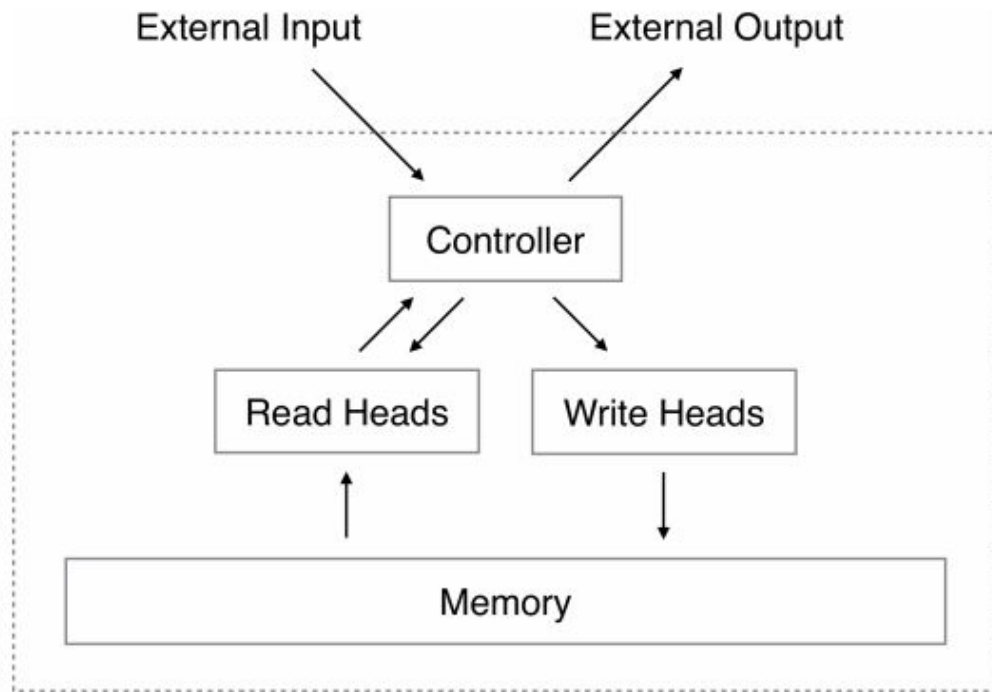


Credit: A. Graves et al.
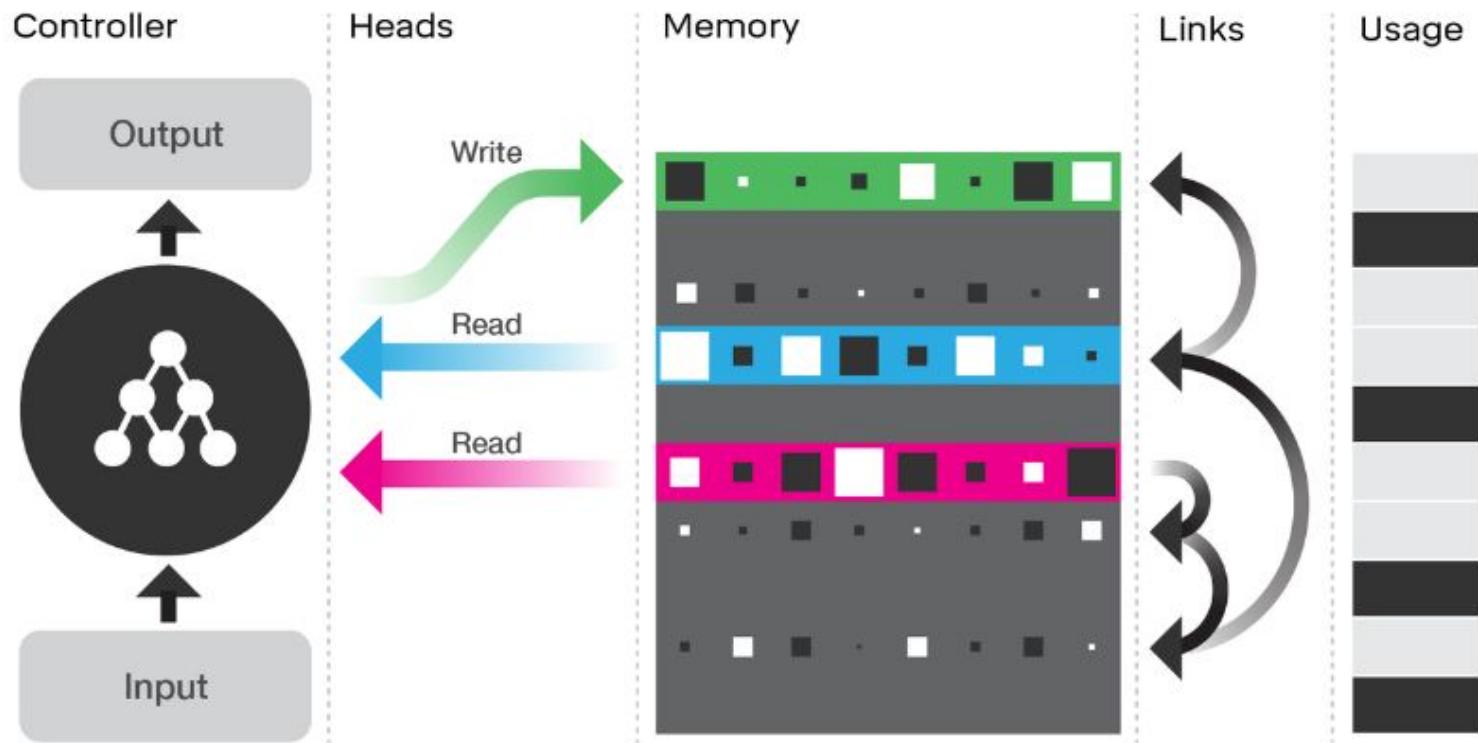
# Illustration of the DNC architecture



Illustration of the DNC architecture. The neural network controller receives external inputs and, based on these, interacts with the memory using read and write operations known as 'heads'. To help the controller navigate the memory, DNC stores 'temporal links' to keep track of the order things were written in, and records the current 'usage' level of each memory location.

# Multi-task: Same network and same core parameters



| Program | Descriptions | Calls |
|---|---|---|
| ADD | Perform multi-digit addition | ADD1, LSHIFT |
| ADD1 | Perform single-digit addition | ACT, CARRY |
| CARRY | Mark a 1 in the carry row one unit left | ACT |
| LSHIFT | Shift a specified pointer one step left | ACT |
| RSHIFT | Shift a specified pointer one step right | ACT |
| ACT | Move a pointer or write to the scratch pad | - |
| BUBBLESORT | Perform bubble sort (ascending order) | BUBBLE, RESET |
| BUBBLE | Perform one sweep of pointers left to right | ACT, BSTEP |
| RESET | Move both pointers all the way left | LSHIFT |
| BSTEP | Conditionally swap and advance pointers | COMPSWAP, RSHIFT |
| COMPSWAP | Conditionally swap two elements | ACT |
| LSHIFT | Shift a specified pointer one step left | ACT |
| RSHIFT | Shift a specified pointer one step right | ACT |
| ACT | Swap two values at pointer locations or move a pointer | - |
| GOTO | Change 3D car pose to match the target | HGOTO, VGOTO |
| HGOTO | Move horizontally to the target angle | LGOTO, RGOTO |
| LGOTO | Move left to match the target angle | ACT |
| RGOTO | Move right to match the target angle | ACT |
| VGOTO | Move vertically to the target elevation | UGOTO, DGOTO |
| UGOTO | Move up to match the target elevation | ACT |
| DGOTO | Move down to match the target elevation | ACT |
| ACT | Move camera $15°$ up, down, left or right | - |
| RJMP | Move all pointers to the rightmost posiiton | RSHIFT |
| MAX | Find maximum element of an array | BUBBLESORT, RJMP |

Credit: N. de Freitas

# Data efficiency: NPI learns to sort with only 8 lessons



**Sorting per-sequence accuracy vs. # training examples**

Legend: ▲ Seq2Seq    ● NPI

X-axis: # Training examples (1, 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024, 2048)

Y-axis: 0, 25, 50, 75, 100

Google DeepMind

# Learning how to learn



parameter updates

optimizer

optimizee

error signal

Credit: N. de Freitas et al.

# References and additional reading

Nando de Freitas' lecture on RNN.

Andrej Karpathy's lecture on RNN.

Deep Learning Book (Chapter 10).

LSTM paper.

Y. Bengio's paper on vanishing gradient problem.

Alex Graves' thesis.

Ilya Sutskever's thesis.

De Freitas' lecture and paper on learning how to learn.

Christopher Olah's educational post on LSTM.

# Thank you!

If you liked this, a series of seminars in deep learning (MLP, CNN, RNN/LSTM, GAN, (Variational) Autoencoders, Reinforcement Deep Learning) is (probably/hopefully) coming next semester (with Sebastiano Vascon).

Deep Learning will be covered and used in the recently formed IEEE Student Branch at DAIS (consider joining it).

We have prepared assignments and thesis proposals on deep learning.