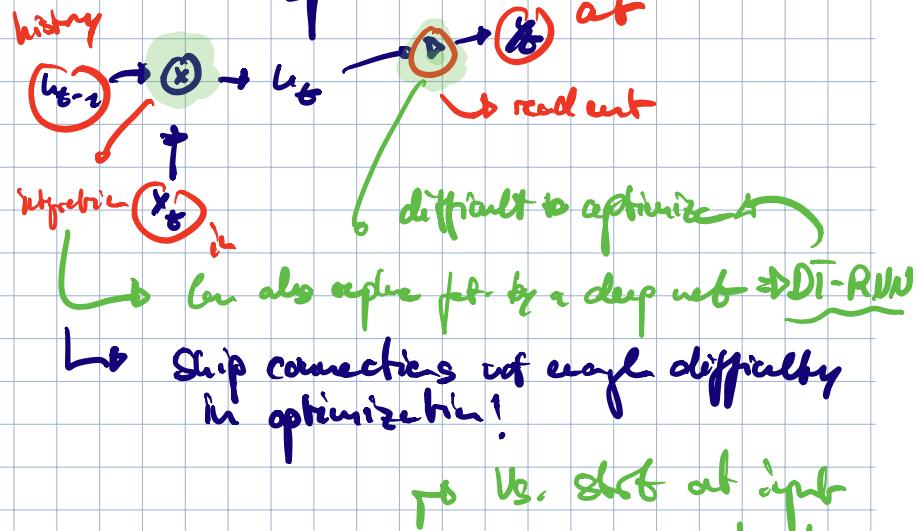
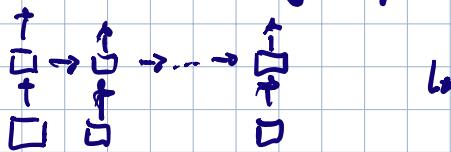


EEML - Day 5 - Bucketize

ULP Task 1 - Razvan + Shyam

- RNN \Rightarrow Parity Complete \rightarrow Expressivity: Can approx. every fct. given finite spaces



- How to make deep?

\rightarrow Stacking:
(for add.
expressivity)

$$\begin{array}{c} \square z_L \\ \uparrow \\ \dots \rightarrow \square z_L \\ \uparrow \\ \dots \rightarrow \square b_L \\ \uparrow \\ \square x_L \end{array}$$

- feedprep: Only have to store vector in every (+) corner complete

$$\frac{\partial L}{\partial x_{a_k}} \in \mathbb{R}^{d_L} \text{ while } \frac{\partial \sigma(\cdot)}{\partial(\cdot)} \in \mathbb{R}^{d_L \times d_L}$$

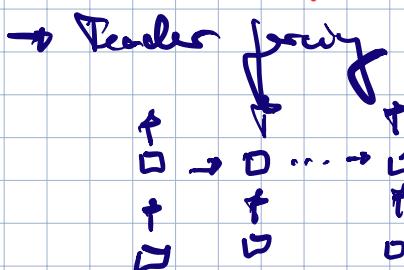
\hookrightarrow Reso: Electrodome architecture give diagonal Jacobian

\hookrightarrow Through time: $\frac{\frac{\partial L(t)}{\partial x(t-h)}}{\frac{\partial L(t-h)}{\partial x(t-h)}} = \prod_{j=h+1}^t \frac{\frac{\partial L(j)}{\partial x(j)}}{\frac{\partial L(j-1)}{\partial x(j-1)}}$ \rightarrow Jacobian product
leads to invertible
matrix prod.
exp. fast

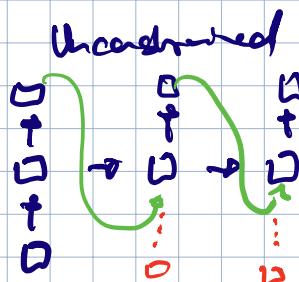
problem: Gradient flow is not enough to solve many

$$\frac{\partial(x_1+x_2+x_3)}{\partial x_3} = 1 \quad \boxed{\text{but}} \quad x_1+x_2+x_3 = 10 \rightarrow \text{Can't reuse } x_3!$$

\hookrightarrow Not all info can be recovered from limited expressivity of hidden state



vs.

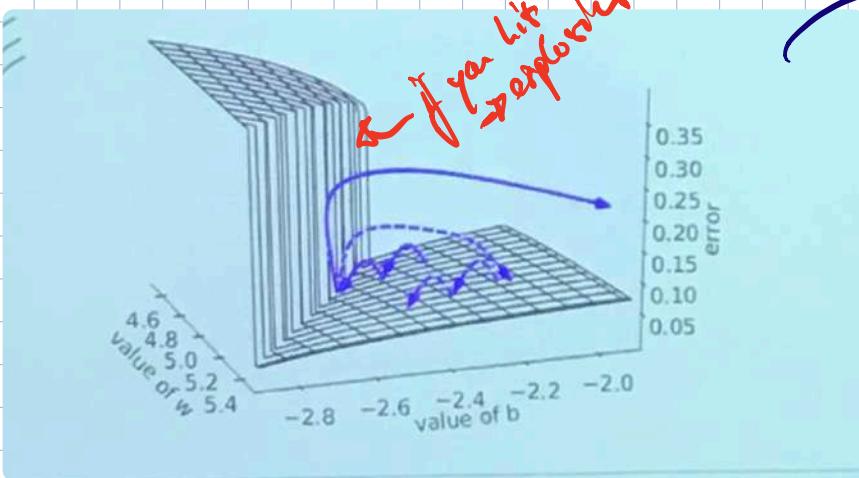


\hookrightarrow looks like Inhibition L.

\hookrightarrow GELERTL: RL \Rightarrow RNN
connection \rightarrow very parallel \rightarrow yield up & freq. trials!

\hookrightarrow can also combine both and then model decides what input to pick!

- Exploding gradients: Broad clipping \Rightarrow ad-hoc solution
 \rightarrow loss before wacky not valley sensitive when gradient / number
 \hookrightarrow More like:



\rightarrow Not solved by 2nd order!
 \hookrightarrow full derivatives explode!

\rightarrow Clipping = Different regions
 \hookrightarrow Clip full gradient
 \hookrightarrow Clip at each layer
 \rightarrow Does not really matter
 \hookrightarrow but you need sensitivity!

- Vanishing gradients: Not slow learning problem,

\rightarrow Components of gradient vanish \Rightarrow makes problem hard
 \hookrightarrow Can't be detected by simply looking at the norm!

$$\hookrightarrow g_i = \frac{\partial C(i)}{\partial x(i)} \rightarrow g_{i-1} = \frac{\partial C(i)}{\partial x(i)} \frac{\partial h(i)}{\partial h(i-1)} + \frac{\partial C(i-1)}{\partial x(i-1)}$$

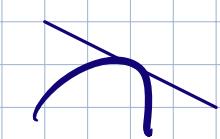
SUM: Problem \rightarrow never have vanishing
 \hookrightarrow parts spread independently! \rightarrow many footprints

\hookrightarrow Also: specifically never explicitly compute $\frac{\partial h(i)}{\partial h(i-1)}$
 \rightarrow exploit element-wise calculations, etc.

- Weight matrix restriction to be orthonormal!

\hookrightarrow Stay on sphere!

\hookrightarrow all eigenvalues = 1



\rightarrow weight happens to explicitly

\rightarrow linearize every term in all directions

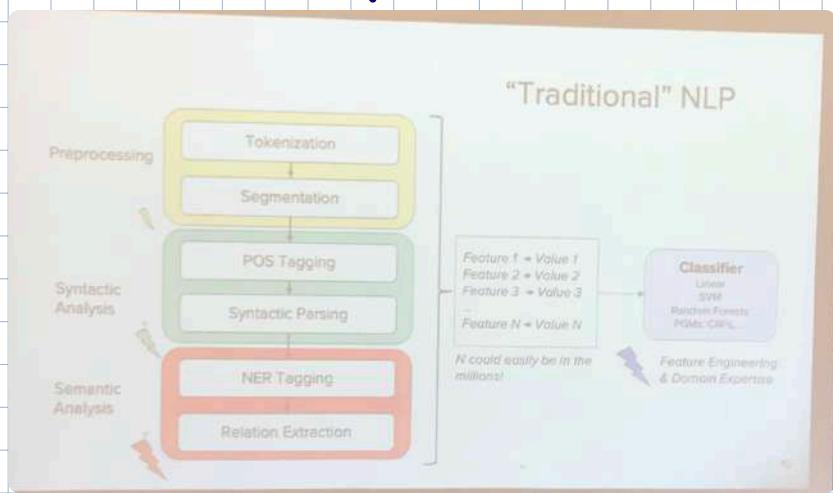
\rightarrow what do we actually learn? directions

\hookrightarrow Saxe et al 2014, Henaff et al 2016, Tishby et al 2016

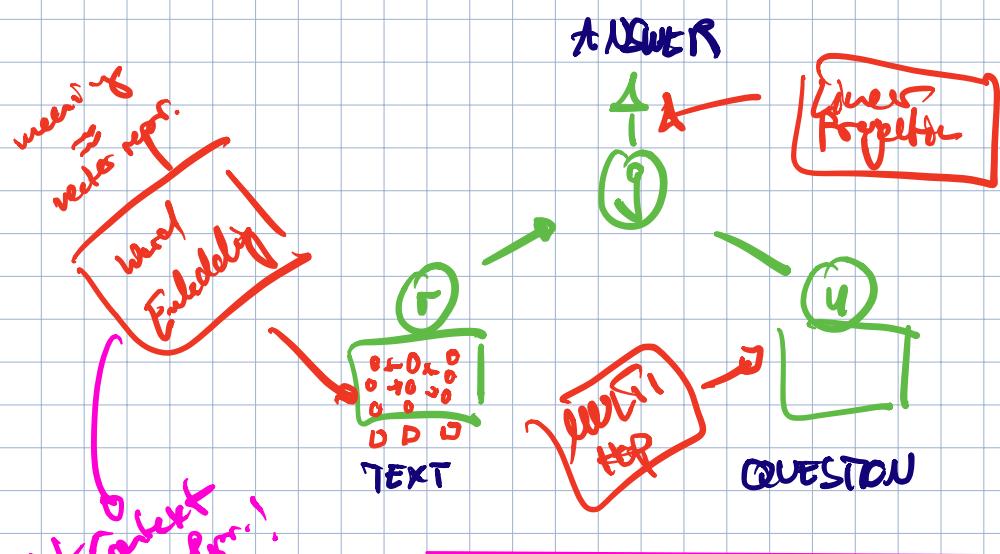
- LSTMs Criticism: Distribution of following the hidden state
 - Non-discriminative gates // How to set gates up?
 - ARNs simplify $\rightarrow h_t = (1 - z) \circ h_{t-1} + z \circ h'$
 - Echo state networks \Rightarrow Learn weights
- Hierarchical approaches: Higher layers work on denser tokens scale
 - Observed LM \Rightarrow Koutnik et al 2014
 - Feeding flags as depth vs **memory loss**
 - \hookrightarrow loss seems to work better!
- Layers:
 - Sutskever et al (2014) \rightarrow Gated
 - Bahdanau et al (2015) \rightarrow Layer
 - \hookrightarrow attention! MP \rightarrow weights \rightarrow aggregate
- WaveNet: Neil Heitzmann - try to use **CNNs** instead of RNNs \rightarrow less steps \rightarrow less memory
 - Easily parallelizable but no longer truly looking at temporal context only act!
 - \hookrightarrow Generalize locality assumption by learning via attention
- Transformer architecture: attention = flexible convolution \rightarrow key
 - Dot product attention: $\star(q, k, v) = \sum_i \frac{e^{q \cdot k_i}}{\sum_j e^{q \cdot k_j}} v_i$ value
 - \hookrightarrow Problem: Softmax \rightarrow amplifies small differences
 - \hookrightarrow Unmodelled \rightarrow Sometimes next word added to update keys
 - \hookrightarrow Same time \rightarrow helps to do credit assignment!
 - \hookrightarrow Solutions: Transformers have separate heads!
 - \hookrightarrow same for key, query, value \rightarrow credit assignment!

NLP Talk 2 - sentence borders - (Machine Reading & Q Answering)

- Machine Reading: Extract representations to answer questions! (from text)
 - before 2014: symbolic approaches → type into SQL query + database
 - after 2014: E-to-E DL
(First RNNs → New Types)
- Challenge: a lot of common sense encoded for us before trying to read!



- Attention Reader Model - Hermann et al 2015



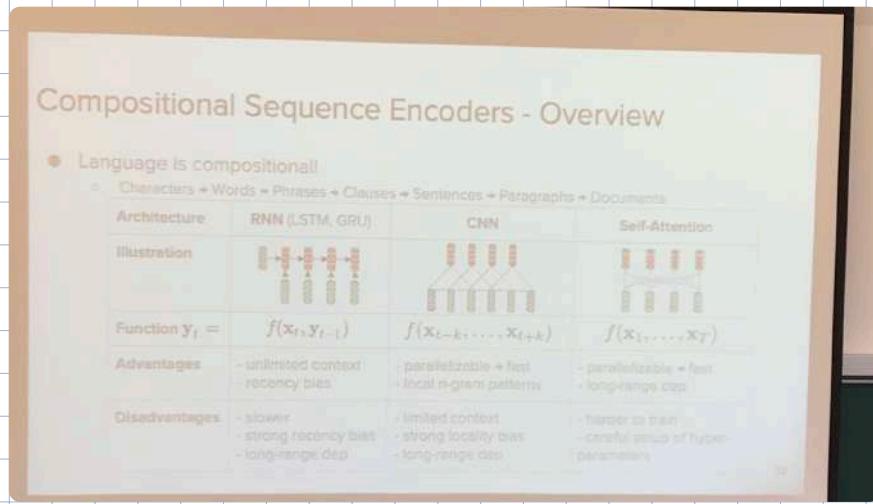
→ Latency \Rightarrow Compositional!
↳ Inductive Biases \Rightarrow with compositionality?

→ Tradeoff \Rightarrow Expressivity vs. Comprehensibility!

→ Bi-Dir. LSTM:

Composable representations
from t -to- t and
 t -to- c

- Self-Selfish \Rightarrow Problems: dependency search, memory costs
 - ↳ Transfer: Vaswani et al 2017 → form graph with weighted edges



written before!

↳ Multi-Headed

⇒ Form of different
Kernels in CNN
without local context!

↳ Can be learned in an
unsupervised fashion!
⇒ PRE-TRAINED!

- Pre-trained Interpretable Embeddings ⇒ ELMo, BERT
 - BERT: randomly mask 15% of tokens in each ⇒ predict!
 - ↳ Learn large transformer to fill blanks ⇒ PROPERTY (G Roberts!)
 - Needs lots of high quality data + lots of compute!
 - ↳ Can also add data in output
 - ↳ Large gap without them
 - ↳ Can only be trained on few computers ⇒ No mobile real-time inference
- Many Webqueries → Multi-hop answer ⇒ sequence matching
 - ↳ Were question explicitly be asked to paragraph?
 - ↳ Go back further (multi-round!)
- answers ⇒ modality distr. over answer options ⇒ linear projection!
 - ↳ Cross-Entropy / L1 loss!
- Query features: paragraph based generation / latent representation / knowledge graph
 - ↳ Das et al., 2013 ⇒ Multi-hop Retriever Reader

EMIL - Day 6 - SparseNet

How to generate stuff and learn representations? - Karl Hechtbräuer

- Generation problem \rightarrow center/focus on info exchange!
- Too many bits in X to model X directly
 - $\Rightarrow p(x) = \prod_i p(x_i | x_{\text{rest}})$ \rightarrow EXPLICIT FACTORIZATION
 - \uparrow Split in pieces (\oplus) simple product pieces independent
 - \downarrow Smaller
 - \rightarrow Reduce modelling capacity!
 - \rightarrow Overfitting

1D \Rightarrow audio + language

- RNN state captures x_t \Rightarrow Worse RNN

8 fine bits



8 coarse bits



\Rightarrow Train very slow!

\hookrightarrow Models optimized for parallel
computable \Rightarrow not separable!

\Rightarrow weights not units!

\rightarrow Sparsifiable \Rightarrow During training \rightarrow residual or add at
batch allows to copy back if gradient in backprop
crosses the threshold

\hookrightarrow Sweet spot of sparsity threshold!

~~SPEED UP!~~
~~Optimize for
Mobile~~

- Subscale Worse RNN: Local dependencies \Rightarrow Global dependencies
 - \hookrightarrow Reshape input tensor with gaps



- Fiber Net: Regularity of end of RNN
 - \hookrightarrow Went to train just \Rightarrow 1d Conv architecture \rightarrow Masked Convolution
- Importance of input features: Granularity vs. Complexity of what to model
 - \hookrightarrow Exploit hierarchy: Characters - Word - Phrases - Words

2D \Rightarrow Vision

- Pixel RNN / Pixel CNN

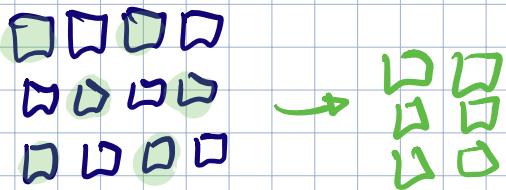
\rightarrow Masked 2D Conv



\rightarrow Squeeze on top
of masked filtering

↳ Model on pixel level!

→ Subscale Pixel Networks: Slicing of image into pieces \Rightarrow conditioned on previously sliced



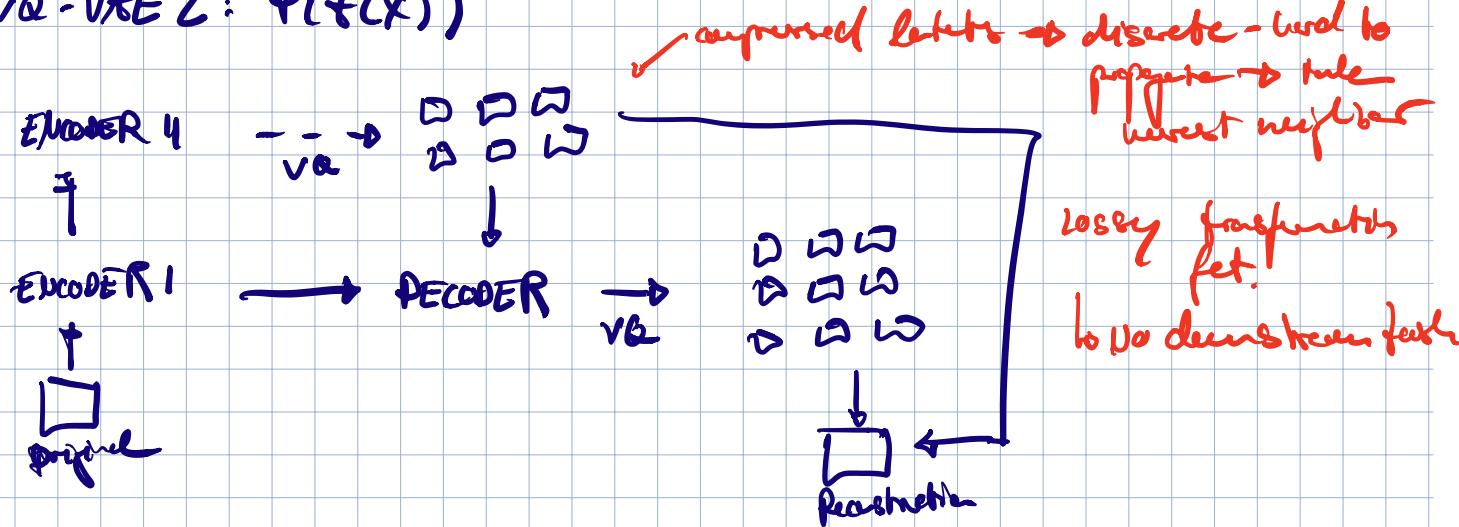
↳ Multidimen. Upscaling: Size + depth \Rightarrow Bits per channel

↳ Was a lot better for images than audio

↳ Working with spectrogram features representation very few channels
 \Rightarrow places looks/sounds fairly random \rightarrow Try models conditionally!

↳ put slices into large encoder-decoder structure

- VQ-VAE 2: $P(f(X))$



→ Latent representations: do not seem to disentangle \Rightarrow Sketch: low-dimensional full sensory representations

3D \Rightarrow Videos

- Video Pixel Networks + Video Transformer

↳ Learned over time
↳ Feed frames and time

New to
generalize!

→ for self-attention
→ predict missing frames

$$f(x; | x_{\leq i}, y_{\geq i})$$

Latent Representations

- Representation Problem: Good fct. $g(X)$ representing X is useful very for classification/recognition

↳ Contrastive Feature Coding

↳ External memory! // need to cluster topics \Rightarrow painful!