Announcements - A4 due tonight - AS out tonight, due next Tuesday Recap Phrase-based MT () Word alignment (src > trg / trg > src, intersect) @ Phrase extraction 3 Decoding search over phrase lattice + LM Seq-to-ser models/Neural MT Today La Connections to RNN LMs Ly Training Ly Problems Attention L. Mechanism Ly Connections to alignment Seq-to-seq models for MT Sutskever et al. Neur IPS Zo14

T t t t č k Je vais le faire whole input sent is source "encoded" into h/c here

svc w/RNN, "decode" try another RNN (diff params) NMT: encode with $P(\mp | \overline{s}) = P(\tau_1 | \overline{s}) P(\tau_2 | \tau_1, \overline{s}) P(\tau_3 | \tau_{12}, \overline{s})$ $P(\tau; | \tau_{1;i-1}, \overline{s}) = \text{Softmax}(W^{(d)}h_{i}^{(d)})$ W^(d): trg Vocob hidden state at ith xhidden size mat step of decoder

Like an RNN LM additionally conditioned on 5

Decoder: feed the previous predicted word into next decoder timestep



Decoding: $t_0 = \angle S7$ while $t_i \neq STOP$ $t_i = \arg\max_{t} P(+(t_{1}, i-1, \overline{s}))$

Looks like LM training Training netrg len trg $w_1 w_2 w_3$ Reall: LM: 1 1 1 RNN CS> w, Uz $loss = \sum_{i=1}^{1} log P(t_i = t_i)$ NMT gradients t, t, t, t, encoder decoder t N P t t t feacher forcing: 5, 52 53 C57 t1 t2 feed in true encoder trained too! target tokens

Teaching the model to predict til the correct this up to This point

What do we really want? Model's predicted sequence f... fn has high BLEU score Requires reinforcement learning

Details Lengths can vary, need to pad MT: 20 chars src try need padding + smart handling

What goes wrong? 1) Repetition Je vais le faire -> I an going to do it going to going to going to ---Wou't happen in PBMI (every word translated as part of one phase) Need a notion of coverage in NMT 2) Unknown words Decoder has a fixed vocab |Vt| PBMT: "copy" rules Pont-de-Buis -> copy to cutput NMT: produce UNK

() Repetition (going to going to) 2) UNKS 3) Poor performance ou long sents BLEU PBMI lo 20 20 40 50 NMT Sent len Bahdanau et al. (2015) (attention) NMT: fixed size h/E vectors don't scale arbitrarily LSTMS can't remember for too long Sutskever: veverse input as a trick to address $S_m S_{m-1} - S_1 \rightarrow t_1 t_2 - t_1 t_2$ the fibst few 5

Attention: Je fais un bureau » I make a desk What if we know translation should be wood-by-wood? Hack I-D-D->D->H-SH-G-J F U b <s> Je fais Un to s, sz Manually feel s; into ith decoder step (along with ti-i) Q: What problems can this fix? (1-3) (+3 / D Less likely to repeat it we see explicit words × 2) Requires changing output layer VB Anchored more to input

I am going to do it i encoder Je vais le faire J V (f.) -Run out of s before end of t? -Ordening is inflexible Attention Target decoder picks where to look in source $\frac{1}{100} + \frac{1}{100} + \frac{1}$ & distribution over source possis

 $m_{(i)}^{scelev} vector$ $m_{(i)}^{i} - (e) vector$ $C_{i} = \sum_{j=1}^{i} \sum_{j=1}^{i} of weighted average$ $j = 1 of h_{i}^{(e)} with a$ as weights $(d) Th_{i}^{(e)}$ $P(t; | 5, t_{(:i-i)}) = softmax \left(W^{(d)} \left[\frac{h_i^{(d)}}{L_i} \right] \right)$ $\overline{X} \quad can also use other fens
to compare <math>\overline{h}^{(d)} / \overline{h}^{(e)}$ besides dot Think about Je fais un bureau $\alpha_{1} = \begin{bmatrix} 0.99 & 0.003 & 0.003 \\ 0.003 & 0.003 \end{bmatrix}$ $Output: softmax \left(W \begin{pmatrix} a \end{pmatrix} & \begin{bmatrix} h & a \\ b & a \\ \hline a & b \end{pmatrix} \right)$ $and ient \qquad \forall Je$ No atta: $\int_{TT} \int_{TT} \int_{TT}$

encoder - Cifer Cifer Je vais le faire I an going to do it Learning mapping 122 443 - Not input at train time ! Je vais le faire -> 1 ZZ YY3 -easier than translation - regular mapping - This is learnable!