"language model" – assigns prob to sequence of text (to even generate text)
 → training data: natural sentences
  → tokenize sentences (one-hot, embeddings).

end token – completes the sequence, learned by LM,
start token – starts sequence, learned by LM.
"finish" rest of sentence – force model to process given segment of sequence.

Conditional LM: instead of $a_0$ (initial state) $= 0$, hidden (activation)
(CNN encoder,                    set $a_0 \leftarrow f(x)$, where (essentially, vector encoding
 RNN decoder).                   $f(x) = CNN$ (input)    of desired content, of sequence.
Training data: (image, words)
This approach is modular: encoder/decoder could be whatever. Ex. encoder is RNN.
  → Sequence to sequence models
   → Typically two separate RNNs, diff weights, trained end-to-end
   → Realistically stacked RNN layers, use LSTM/GRU cells. diff. length sequences.
   → utilizations: translation, summarization, Q/A, Text-to-code.
   → Typically we reverse order of tokens to encoder (due to how memory is processed)

Decoding likely sequence: maximize product of all sequence probabilities.
Each seq output: $p(y_{i,t} | x_{i, 1:T_i}, y_{i, 0:t-1}) \leftarrow$ conditioned on everything before it.
So $p(y_{i,1:T_y} | x_{i,1:T}) = \prod_{t=1}^{T_y} p(y_{i,t} | x_{i, 1:T_i}, y_{i, 0:t-1}) \leftarrow$ product of conditional probs
                                                                          (chain rule)
# decodings: for M words $M^T$ seqs for T length. Turns into search problem! ←
  → expensive to find optimal seq. using exact method.
Approximate search to find largest prod: avoid super low probs. Target top k probs @ init.
 → "Beam Search" – store k best sequences so far, update each of them.
  at each t,                                          ↙ k of these
   1. for each hypothesis $y_{1:t-1, i}$ that we are tracking:
    ² find top k tokens $y_{t, i, 1}, ..., y_{t, i, k}$
   2. sort resulting $k^2$ length t sequences by total log probs
   3. keep top k
   4. advance each hypothesis to time t+1 (forward)
 → If top sequences ends abruptly, save it. "remove" it. move on.       ← predetermined.
 → continue until $t > T$ or # of seqs that end in $<EOS> > N$.
 → perhaps divide total log prob by T to generate "score" for seq (longer seq → more penalty).

"Bottleneck" problem → for decoder, all information stored in potentially vast encoder → one layer
  → instead, while decoding, "peek" at source sentence.                          of activations
  → for each x layer in encoder, generate key vector thru learned function. (what info is present)
  → for each y layer in decoder, generate query vector thru function (what info we are
  → compare query to key to find closest one.                                    looking for)
                                                        ← attention score
Mathematically: $k_t = k(e_t)$, $q_i = q(d_i)$, $s_{t,i} = k_t \cdot q_i$ for encoder step to decoder step.
Then, to get max $s_{t,i}$ use softmax (argmax is ideal but not differentiable)
$\alpha_{t,i} = softmax(s_{:,i})$, $\alpha_{t,i} = \frac{exp(s_{t,i})}{\sum_{t'} exp(s_{t',i})}$ ← softmax      ★ Network is trying to "pay
                                                                                       attention" to most relevant
"send" $a_i = \sum_t \alpha_{t,i} e_t \leftarrow$ approx. to "best" $e_t$ for given $d_i$     part of input.
(i.e. $\hat{y}_i = f(d_i, a_i)$, pass into RNN layer, pass into next decoder step $\left( \bar{d}_i = \begin{bmatrix} d_{i-1} \\ a_{i-1} \\ x_i \end{bmatrix} \right)$
  (use for readout)                     (concatenate to hidden state)
Attention Variants
 → $k_t = e_t$, $q_i = d_i$, so $s_{t,i} = e_t \cdot d_i$        learn this   → learned-value encoding:
                                                              ↓              $a_i = \sum_t \alpha_{t,i} v(e_t)$
 → $k_t = W_k e_t$, $q_i = W_q d_i$, so $s_{t,i} = e_t^T W_k^T W_q d_i = e_t^T W_s d_i$      ↖ learn this