# Generation

Vsevolod Dyomkin
prj-nlp 2018-05-16

# Natural Language Generation (NLG)

* general-purpose
* special-purpose

* word-level
* text-level
* book-level

Applications:
* data-to-text
* simplification
* summarization
* paraphrasing
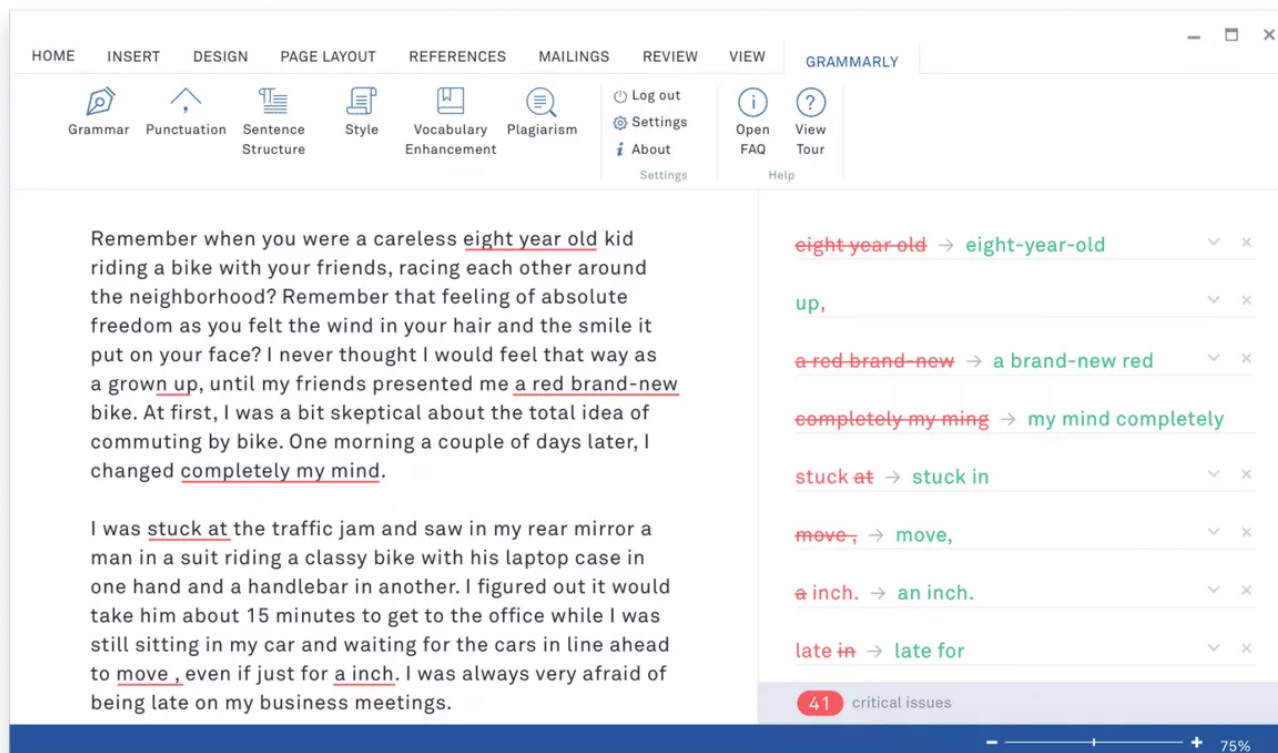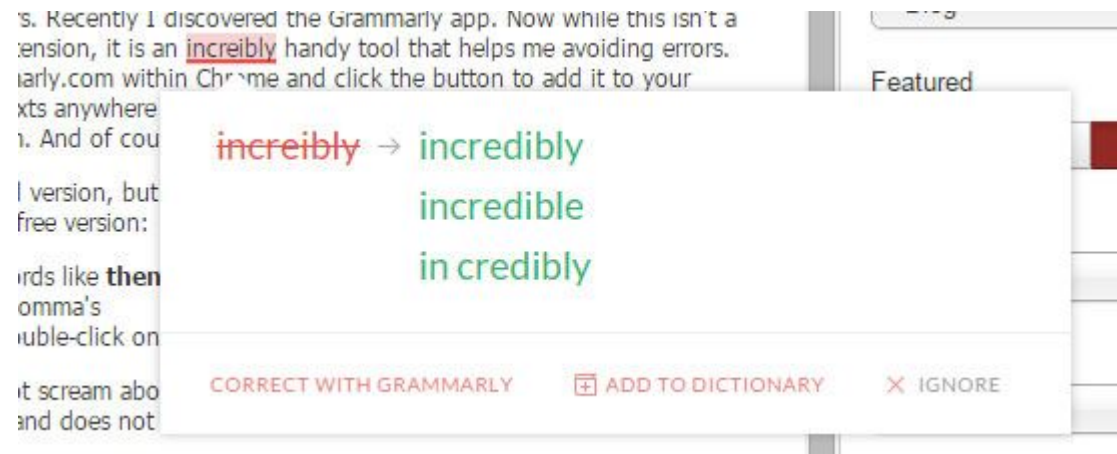* dialogue
* computer-generated verse/poetry

* MT
* GEC
* QA

# Levels of NLG

Level 1: Simple Fill-in-the-Blank Systems
Level 2: Script/Rule-based Systems
Level 3: Word-Level Grammatical Functions
Level 4: Dynamically Creating Sentences
Level 5: Dynamically Creating Documents

https://ehudreiter.com/2016/12/18/nlg-vs-templates/

# Example: Grammarly

* word choice
* phrase rewriting
* … text rewriting

s. Recently I discovered the Grammarly app. Now while this isn't a
ension, it is an <u>increibly</u> handy tool that helps me avoiding errors.
arly.com within Chrome and click the button to add it to your
xts anywhere
. And of cou

~~increibly~~ → incredibly
incredible
in credibly

I version, but
free version:

rds like **then**
omma's
uble-click on

t scream abo
and does not

Featured

CORRECT WITH GRAMMARLY      ADD TO DICTIONARY      ✕ IGNORE

HOME   INSERT   DESIGN   PAGE LAYOUT   REFERENCES   MAILINGS   REVIEW   VIEW   GRAMMARLY

Grammar   Punctuation   Sentence Structure   Style   Vocabulary Enhancement   Plagiarism

⏻ Log out
⚙ Settings
ℹ About
Settings

ⓘ Open FAQ
? View Tour
Help

Remember when you were a careless <u>eight year old</u> kid riding a bike with your friends, racing each other around the neighborhood? Remember that feeling of absolute freedom as you felt the wind in your hair and the smile it put on your face? I never thought I would feel that way as a <u>grown up</u>, until my friends presented me <u>a red brand-new</u> bike. At first, I was a bit skeptical about the total idea of commuting by bike. One morning a couple of days later, I changed <u>completely my mind</u>.

I was <u>stuck at</u> the traffic jam and saw in my rear mirror a man in a suit riding a classy bike with his laptop case in one hand and a handlebar in another. I figured out it would take him about 15 minutes to get to the office while I was still sitting in my car and waiting for the cars in line ahead to <u>move ,</u> even if just for <u>a inch</u>. I was always very afraid of being late on my business meetings.

~~eight year old~~ → eight-year-old                    ∨ ✕

up,                                                    ∨ ✕

~~a red brand-new~~ → a brand-new red                  ∨ ✕

~~completely my ming~~ → my mind completely            ∨ ✕

stuck ~~at~~ → stuck in                                 ∨ ✕

~~move ,~~ → move,                                      ∨ ✕

~~a inch.~~ → an inch.                                  ∨ ✕

late ~~in~~ → late for                                  ∨ ✕

(41)  critical issues

−  |  +  75%

# Example: data-to-text

* weather
* sports
* stocks
* news


https://www.aclweb.org/anthology/W18-6504

# Example: chatbots

verloop

# Example: books



**Darius Kazemi**
@tinysubversions

Hey, who wants to join me in NaNoGenMo: spend the month writing code that generates a 50k word novel, share the novel & the code at the end

♡ 179   7:00 PM - Nov 1, 2013

ℹ

💬 148 people are talking about this   >

Titles:
* Webster's Slovak — English Thesaurus Dictionary for $28.95
* The 2007-2012 World Outlook for Wood Toilet Seats for $795
* The World Market for Rubber Sheath Contraceptives (Condoms): A 2007 Global Trade Perspective for $325
* Ellis-van Creveld Syndrome — A Bibliography and Dictionary for Physicians, Patients, and Genome Researchers for $28.95
* Webster's English to Haitian Creole Crossword Puzzles: Level 1 For $14.95

https://singularityhub.com/2012/12/13/patented-book-writing-system-lets-one-professor-create-hundreds-of-thousands-of-amazon-books-and-counting/

# Example: abstractive summarization

| | |
|---|---|
| **Article** | novell inc. chief executive officer eric schmidt has been named chairman of the internet search-engine company google . |
| **Human summary** | novell ceo named google chairman |
| **Textsum** | novell chief executive named to **head** internet company |

https://rare-technologies.com/text-summarization-in-python-extractive-vs-abstractive-techniques-revisited/

# NLG Evaluation

Need to capture quality & diverscity

(Best) Real-World Task-Based (Extrinsic)
(Good) Laboratory Task-Based or Real-World Human Ratings
(OK) Laboratory Human Ratings
(Worst) Metrics

https://ehudreiter.com/2017/01/19/types-of-nlg-evaluation/

# Metrics

* BLEU
* ROUGE
* METEOR
* Perplexity

# NLG Evaluation: HUSE

Combine human evaluation & perfplexity



https://arxiv.org/pdf/1904.02792.pdf

# Classic Approach to NLG

1) Content determination
2) Document structuring
3) Aggregation
4) Lexical choice
5) Referring expression generation
6) Realization

# Hybrid Approaches

* overgenerate than select
https://aclanthology.info/pdf/P/P98-1116.pdf
(an example using AMR)

* ML choosers embedded in a rule-based framework
https://aclanthology.info/pdf/J/J17-1001.pdf

https://ehudreiter.com/2017/10/16/machine-learning-and-rules/

# DL Approaches

* a plain RNN
* variational autoencoders
* seq2seq
* transformers

Also, maybe:
* GANs
* deep re-inforcement learning

# VAEs

A generation model "framework":
- encoder
- hidden state
- decoder



(a) VAE training graph using a dilated CNN decoder.



(b) Digram of dilated CNN decoder.

# Language Modelling Task

Question: what is the probability of a sequence of words (sentence/paragraph/text)?

And why do we need it?

For the sentence

the dog barks STOP

we would have

$$p(\text{the dog barks STOP}) = q(\text{the}|*, *)$$
$$\times q(\text{dog}|*, \text{the})$$
$$\times q(\text{barks}|\text{the, dog})$$
$$\times q(\text{STOP}|\text{dog, barks})$$

# LM Applications

* Word choice, predictive typing
* NLG
* Statistical machine translation
* Spelling & grammatical error correction
* OCR, ASR, code breaking, paleolinguistics etc.
* transfer learning

# Ngram LM

Apply Markov assumption to the word sequence.

If n=3 (trigrams):
$$P(S) = P(w0) * P(w1|w0) * P(w2|w0\ w1)$$
$$* P(w3|\cancel{w0}\ w1\ w2) * P(w4|\cancel{w0\ w1}\ w2\ w3)$$

According to the chain rule:

$$P(w2|w0\ w1) = P(w0\ w1\ w2) / P(w0\ w1)$$

We can use MLE

# Ngrams Estimation



The most frequent Wikipedia bigrams beginning with 'red'

# Ngrams' Problems

* Need big corpus for MLE
* Number of ngrams ~ $O(e^n)$ (n-ngram rank)
* Sparsity (problem of UNKs):

$$P(S) = P(w0) * P(w1|w0) * P(w2|w0\ w1)$$
$$* P(w3|w1\ w2) * P(w4|w2\ w3)$$

If some of w0-w4 are UNK P(S) = 0!

# Ngrams Smoothing

* Laplace smoothing

# Ngrams Smoothing

* Laplace smoothing
* Naive +1 smoothing

# Ngrams Smoothing

* Laplace smoothing
* Naive +1 smoothing
* Good-Turing smoothing, Katz smoothing
* Knesser-Ney smoothing:
    a discounting interpolation
    (using lower-order ngrams)

$$P_{KN}(w_i \mid w_{i-1}) = \frac{\max(c(w_{i-1}w_i) - \delta, 0)}{\sum_{w'} c(w_{i-1}w')} + \lambda \frac{|\{w_{i-1} : c(w_{i-1}, w_i) > 0\}|}{|\{w_{j-1} : c(w_{j-1}, w_j) > 0\}|}$$

$$\lambda(w_{i-1}) = \frac{\delta}{c(w_{i-1})} |\{w' : c(w_{i-1}, w') > 0\}|$$

# Ngrams Implementation

* cut-off
* efficient storage (binary trees, perfect hash-tables, …)
* quantization
* efficient estimation (MapReduce)

LM Software:
* BerkeleyLM
* KenLM

https://kheafield.com/papers/stanford/crawl_paper.pdf

# LMs Evaluation

Intrinsic evaluation -
perplexity (a measure of surprise /per word):

$$2^{H(p)} = 2^{-\sum_x p(x) \log_2 p(x)}$$

$$PP(s_1, s_2, ...) = \sqrt[(\sum_i |s_i|)]{\frac{1}{\prod_i p(s_i)}}$$

A corpus-based measure. Current corpus — 1B
word benchmark (http://arxiv.org/abs/1312.3005)

Extrinsic evaluation also necessary

# SOTA Perplexity

| Model | Test Perplexity |
|---|---|
| Sigmoid-RNN-2048 (Ji et al., 2015a) | 68.3 |
| Interpolated KN 5-gram, 1.1B n-grams (Chelba et al., 2013) | 67.6 |
| Sparse Non-Negative Matrix LM (Shazeer et al., 2015) | 52.9 |
| RNN-1024 + MaxEnt 9-gram features (Chelba et al., 2013) | 51.3 |
| LSTM-512-512 | 54.1 |
| LSTM-1024-512 | 48.2 |
| LSTM-2048-512 | 43.7 |
| LSTM-8192-2048 (No Dropout) | 37.9 |
| LSTM-8192-2048 (50% Dropout) | 32.2 |
| 2-Layer LSTM-8192-1024 (Big LSTM) | 30.6 |
| Big LSTM+CNN Inputs | **30.0** |
| Big LSTM+CNN Inputs + CNN Softmax | 39.8 |
| Big LSTM+CNN Inputs + CNN Softmax + 128-dim correction | 35.8 |
| Big LSTM+CNN Inputs + Char LSTM predictions | 47.9 |

https://arxiv.org/pdf/1602.02410.pdf

# Character LM

What if we use characters instead of words (for ngrams or as input to the NN)?

... "The unreasonable effectiveness of Character-level Language Models"
http://nbviewer.jupyter.org/gist/yoavg/d76121dfde2618422139

For ngram-based models, as number of tokens is small, order may be quite large (10-20-100?)

Pro: no need for smoothing
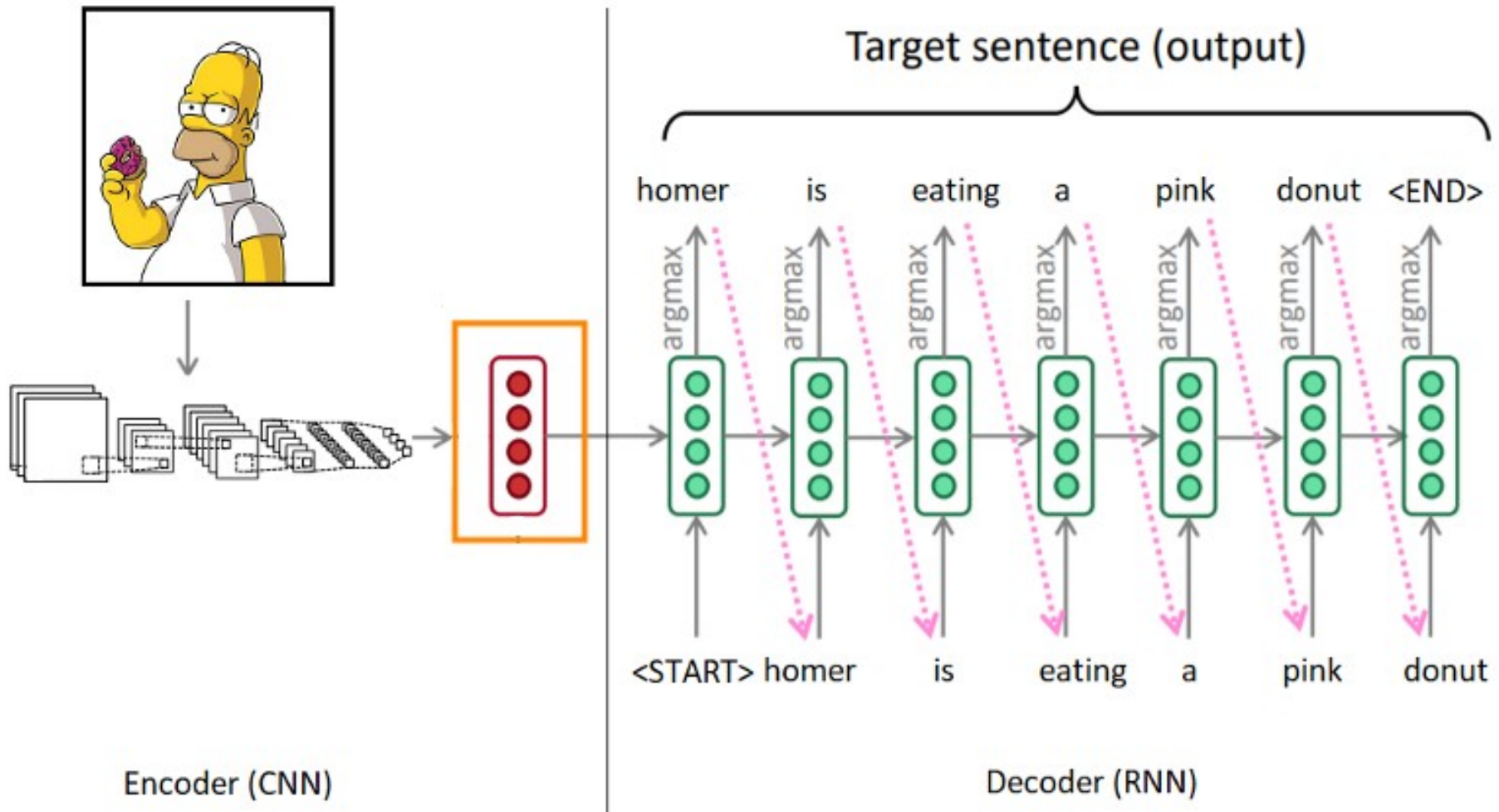Con: no notion of tokens

# Neural CharLM

# Neural LM



http://www.jmlr.org/papers/volume3/bengio03a/bengio03a.pdf

# seq2seq

# seq2seq variants: image captioning

# seq2seq variants: guided generation

# Attention

# LMs Recap

LMs may be used both in classification and generation tasks:
* in classification they can be combined with a domain model
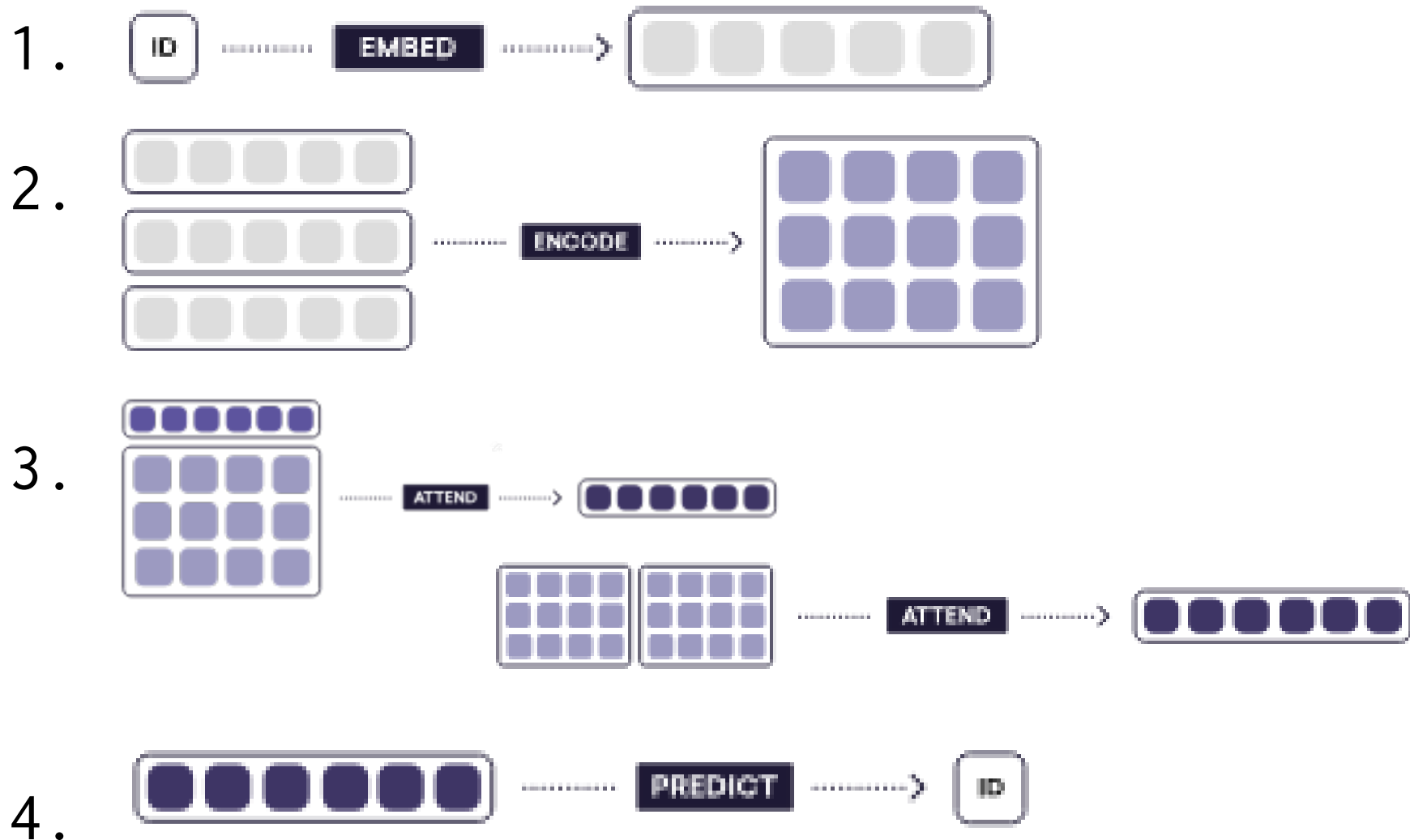* in generation: sample from the model or re-rank other model's output

Main approaches:
* charLMs
* smoothed ngrams
* neural language models (SOTA)
* but other variants are also possible (grammars, topic models…)

# The "DL Formula"

Embed, encode, attend, predict

# NLG Recap

* NLG - the pinnacle of NLP
* Allows for many approaches.
  A good area to utilize DL strong points.
* But evaluation is complicated
  (+ lack of quality resources)

# Read More

NLG:

https://ehudreiter.com
https://arxiv.org/pdf/1509.00685.pdf
https://aclweb.org/anthology/J/J12/J12-1006.pdf
https://www.youtube.com/watch?v=9zKuYvjFFS8

LMs:

http://www.dhgarrette.com/nlpclass/notes/ngrams.pdf
http://www.foldl.me/2014/kneser-ney-smoothing/

NNs:

http://www.wildml.com/2016/01/attention-and-memory-in-deep-learning-and-nlp/
https://medium.com/@yoav.goldberg/an-adversarial-review-of-adversarial-generation-of-natural-language-409ac3378bd7
https://medium.com/@hyponymous/paper-summary-neural-machine-translation-by-jointly-learning-to-align-and-translate-84970177e08c
http://ofir.io/Neural-Language-Modeling-From-Scratch/
https://slides.com/oleksiysyvokon/lm-advances
https://medium.com/@adityathiruvengadam/transformer-architecture-attention-is-all-you-need-aeccd9f50d09