# Language Model Evaluation

## Splitting the Data

We will now discuss the train/val/test splits and perplexity.

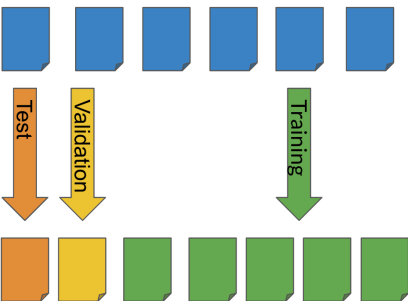**Train/Val/Test splits**

Smaller Corpora:

- 80% train
- 10% val
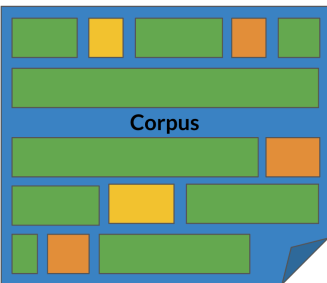- 10% test

Larger Corpora:

- 98% train
- 1% val
- 1% test

There are two main methods for splitting the data:

- Continuous text
- Random short sequences

## Perplexity

Perplexity is used to tell us whether a set of sentences look like they were written by humans rather than by a simple program choosing words at random. A text that is written by humans is more likely to have lower perplexity, where a text generated by random word choice would have a higher perplexity.

Concretely, here are the formulas to calculate perplexity.

$$PP(W) = P(s1, s2, \ldots, sm)^{-\frac{1}{m}}$$

$$PP(W) = \sqrt[m]{\prod_{i=1}^{m} \prod_{j=1}^{|si|} \frac{1}{P(w_j^{(i)} | w_{j-1}^{(i)})}}$$

$w_j^{(i)} \rightarrow$ j corresponds to the jth word in the ith sentence. If you were to concatenate all the sentences then $w_i$ is the ith word in the test set. To compute the log perplexity, you go from

$$PP(W) = \sqrt[m]{\prod_{i=1}^{m} \frac{1}{P(w_i | w_{i-1})}}$$

To

$$\log PP(W) = -\frac{1}{m} \sum_{i=1}^{m} \log 2 \left( P(w_i | w_{i-1}) \right)$$