# A STUDY ON LINGUISTIC COMPLEXITY

Robin Malhotra

# INSPIRATION

- The Hemingway App (http://www.hemingwayapp.com)

# CURRENT STATE

← → C 🔒 jsbeautifier.org

Beautify, unpack or deobfuscate JavaScript and HTML, make JSON/JSONP readable, etc.

All of the source code is completely free and open, available on GitHub under MIT licence,
and we have a command-line version, python library and a node package as well.

Indent with 4 spaces
Allow 5 newlines between tokens
Wrap lines near 120 characters
Braces with control statement

HTML <style>, <script> formatting:
Add one indent level

☐ End script and style with newline?
☐ Support e4x/jsx syntax
☐ Use comma-first list style?
☑ Detect packers and obfuscators?
☐ Keep array indentation?
☐ Break lines on chained methods?
☑ Space before conditional: "if(x)" / "if (x)"
☐ Unescape printable chars encoded as \xNN or \uNNNN'
☐ Use JSLint-happy formatting tweaks?
☐ Indent <head> and <body> sections?
Use a simple textarea for code input?

**Beautify JavaScript or HTML** (ctrl-enter)

```
396        orig: _0xb107[220],
397        replace: _0xb107[221]
398    }, {
399        orig: _0xb107[222],
400        replace: _0xb107[223]
401    }, {
402        orig: _0xb107[224],
403        replace: _0xb107[225]
404    }, {
405        orig: _0xb107[226],
406        replace: _0xb107[226]
407    }, {
408        orig: _0xb107[227],
409        replace: _0xb107[228]
410    }, {
411        orig: _0xb107[229],
412        replace: _0xb107[229]
413    }, {
414        orig: _0xb107[230],
415        replace: _0xb107[231]
416    }, {
417        orig: _0xb107[232],
418        replace: _0xb107[233]
419    }, {
420        orig: _0xb107[234],
421        replace: _0xb107[234]
422    }, {
423        orig: _0xb107[235],
424        replace: _0xb107[235]
425    }, {
426        orig: _0xb107[236],
427        replace: _0xb107[236]
428    }, {
429        orig: _0xb107[237],
430        replace: _0xb107[237]
431    }, {
432        orig: _0xb107[238],
433        replace: _0xb107[239]
```

Beautify JavaScript or HTML (ctrl-enter)

# THE NOTION OF READABILITY

- What makes text readable?

- Is there a metric to capture readability?

# APPROACHES
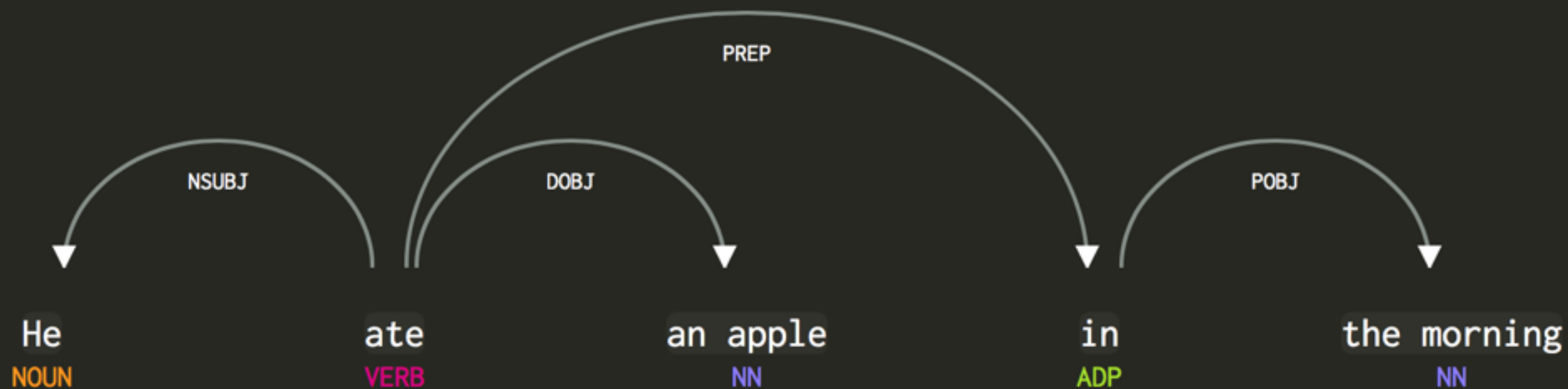
- Word Frequency: unigram probability for a word to appear.

- More frequent words => difficult to read

# APPROACHES

- Structural : Words with deeper, 'more complicated' trees tend to be harder to read.

  Ex: "He, in the morning, ate an apple" (p=0.04) is more complex than
  "He ate an apple in the morning" (p=0.019)

# SPLACY.IO

## How does spaCy compare to NLTK?

**SPACY**

- ■ Over 400 times faster
- ■ State-of-the-art accuracy
- ■ Tokenizer maintains alignment
- ■ Powerful, concise API
- ■ Integrated word vectors
- ■ English only (at present)

**NLTK**

- ■ Slow
- ■ Low accuracy
- ■ Tokens do not align to original string
- ■ Models return lists of strings
- ■ No word vector support
- ■ Multiple languages

Figure 1

Histogram of NormalisedDependancyLengths : $\mu = 2.541$, $\sigma = 0.578$