

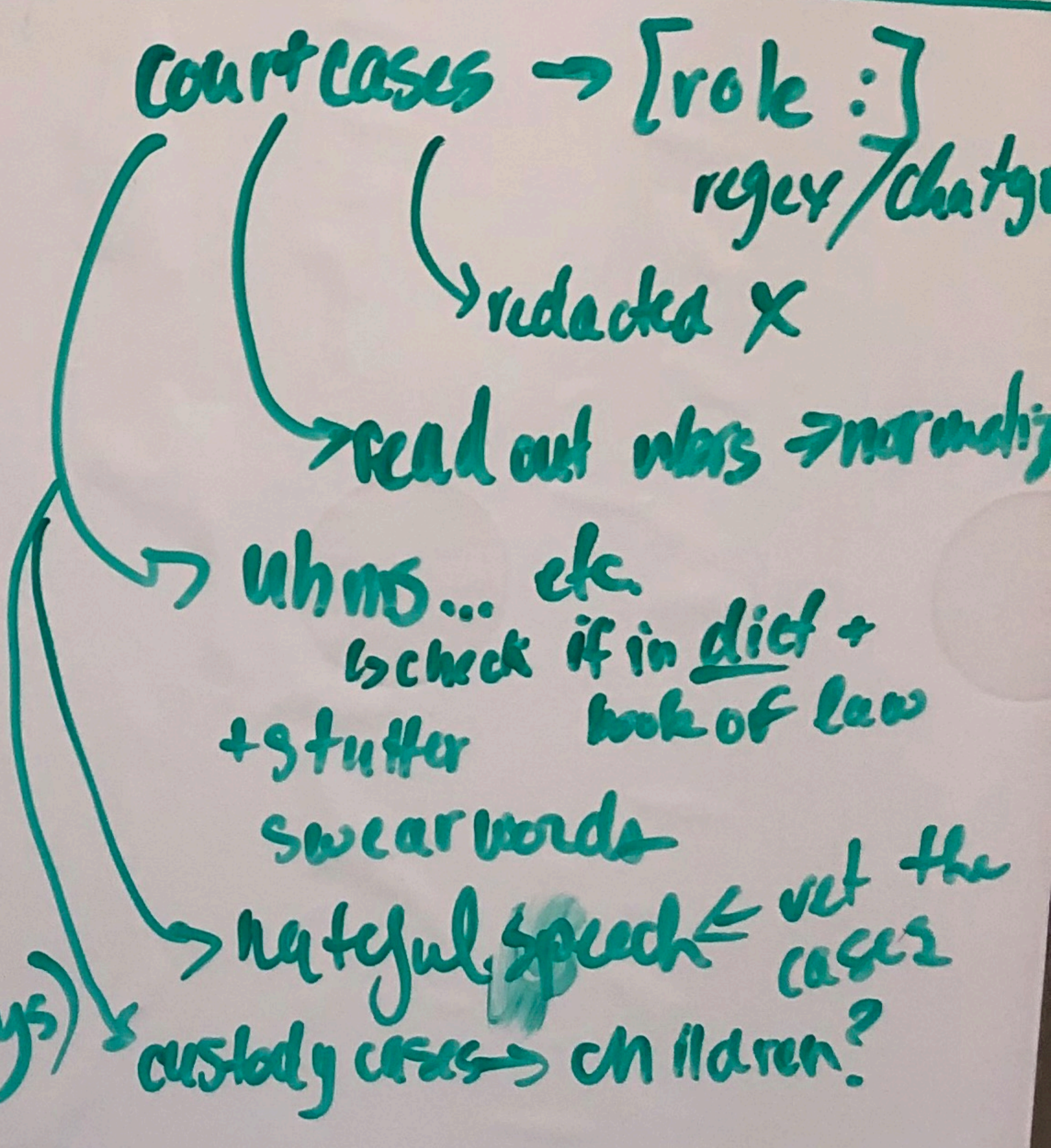
DOMAIN - source

- reddit (typos, Native S. vs NON-N, upvotes)
- transcribed press conferences (EU)
- copyright free poems / lyrics

- Wikipedia (Scientific Vs colloquial, topic moderated)
- old movie scripts
- magazines
- children's books (might get nonsense w/ # Fairytales)
- personal emails
- journal entries from 2018
- Court cases (Shakespeare, King James Bible)
- podcasts (w/ copyright)

CLEAN

- Wiki (formulae + greek letters, Currency, symbols, read out or take out, foreign names/h-pronounce w/feeling)
- children's books + fairytales (Names, Countries, nonsense, feed to dict. + pop)
- Shakespeare (w/CMV dict poetry, plays)



RICHNESS

- # unique phones (not already found)
- normalize by # diphones
- weight by diphone freq
- consider sentence position / word position
- prosody - sentence type
- separate lists of diphones for pop diphones as encounter



Step 1: Source - Out of copyright books, travel blogs,

Speechzone, open access books, Wikipedia,

AI-generated text, Common Crawl, twitter, Stack overflow

Parliament Preceding

Step 2: Clean

(reusability) & filter away too long sentences

Short

hashtags & emojis
numbers
dates
acronym...

text-normalisation

programming -> formatted text

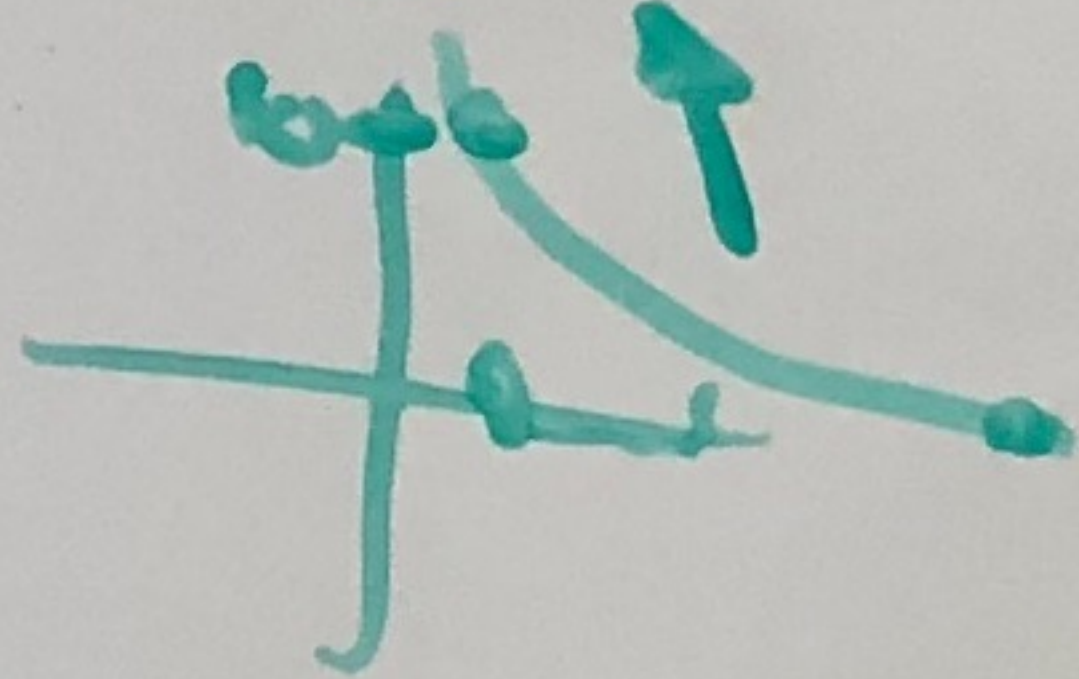
non-Eng &

(synthesiser's frontend)

↓ further (selection)

diversify probably by

sentence types & phrase positions



Step 3:

- diphones (weighted according to English freq. & how many we have)
- length / # of words
- question / types of prosodic inf.
- neighboring diphones
- stress / unstressed
- pitch / fo low / high

Full coverage
dynamic cost

the kat sat

[t a s k a e t s a e t]

zoo z k-ae-stressed

o, + 'o + z

STEP 1.

+ non-comm License

SOURCE: • TED TALKS

- NOW corpus + free + recent (2010) + big
 - Amazon Q&A data + just cite - informal - names ...
- + readability
+ quality, transcription
- narratives

STEP 2.

- F.E.

- Out-of-vocab words, names
- Remove too long/short
- clean away markups, timestamps
- Get rid of question marks, etc. (exclamation)

STEP 3.

• weighted sum of interesting features

- phoneme (type) coverage

$\sum W_i d_i$

- stress/unstress

- position in syllable

- position of punctuation

- phonetic context

- compare w/
existing data
higher - if it
adds variety

Source Text:

Tweets

- Short
- Text Norm is Awful
- Prosody ↑

Academic Papers

- Domain Specific
- Prosody ↓
- Hard to Access

"General Science"

Textbooks

- Variety of Domain

EAS Systems

- Short

Myths

~~Entities~~

Exp-to-Speed

Clean Text

Year & Dates
 Figure Numbering

Named Entities

Equations

Individual Blocks

~~Pronunciation~~

Punctuation

Figure Caption

Measurement Units

Abbreviation

Score [0-1] (Use python G2P) The cat sat.

Num. of phones/diphones

$[v_1]$
 $[v_2]$
 ↑

[variety, length]

T < 15 phonemes

SK @ ə

↓
 kæt sæt sk = 10
 ↑
 = 20

10 + 10

PL: [~~~~~] # 100

School website
Government website: gov.uk

public dataset: [] speech. LibriTTS

~~English~~ vedna.c

Step 2 clean

HTML

1) script to find text

2) text normalization: abbreviation? capitalized?

· Acronyms/initialisms

Spelled-out

· Decide line length/segment/Number

Step 3

· Use Arctic as baseline

· Count diphone occurrences

· Use that as weights

(less frequent → higher score)

· Give "?" / "!" more points...

· Normalize by total # of diphones

twitter corpus (?)

reddit

wikipedia

Spotify (podcasts)

interviews

+ data
conversational
modern

- NSW
(too many)

- toxic

TTS
standard language

- not conversational

- copyright

- topic domain

REMOVE "#s" "@s" ^{links}

• tweets already split data into readable chunks

• rank by length

expand? "jk" → "just kidding"
"lol" ??? → "but different."
ES: "jkhjkhj...."

"AI", "NLP"?

• multilingual tweets?

• typos

⑤ - Perplexity - range

BUT we want to keep rare sounds...

② { sentence length }
{ + character length }
{ + syllables } → ③ Prioritise rare diphones

① - how much 'cleaning' has been done

④ - Prioritise questions for prosody range

source of txt:

	COCA	BNC
readability	21st but too long 1.5/2	outdated (20th) 1/2
size	1 billion 1/2	100 million 1/2
copyright	✓ =	✓ =
NSWs	✓	✓

Clean the text:

- ~~Shorten sentences?~~
- Filter out hard to pronounce words & misspellings.
- Include sentences with uncommon and NSWs.
- Filter according length (10)
 - ↳ different domains.
- sentences with: numbers, homographs, etc.
- Diverse punctuation contexts

Richness Measure

- frequency of words → proportion to real life
MLE
- Depth of syntactic tree (variation) (n° nodes + level)
- Diphones, capturing contextual effects
- No of breaks
- N° of diphones → $\frac{c(\text{diff-diph})}{c(\text{diph.})}$

Machine learning methods (Random Forest → importance of features)

$$f(x) = G(\dots)$$

Sources of Text

1. Parliament Proceedings. [Europarl] (lots of NSWs, natural, ✓)
2. Wikipedia (lots of NSWs, lots of processing, ✓)
3. Dictionary (unnatural)
4. Twitter (lots of NSWs, natural, ↑ preprocessing, ↓)
5. Encyclopaedia (domain ✓, natural, ↓ NSWs, ↑ variety)
6. News (natural speech, domain ✓, NSWs)
7. Call centre scripts
8. Committee meeting minutes

1. Scrape, normalise, reformat → script prosodic weirdness
2. ✓
3. format, preprocessing, TEXT ≠ SPEECH
4. Tools available? time cost
5. is it real? ~~ASR?~~ ~~reliability?~~

'What is our sentence'

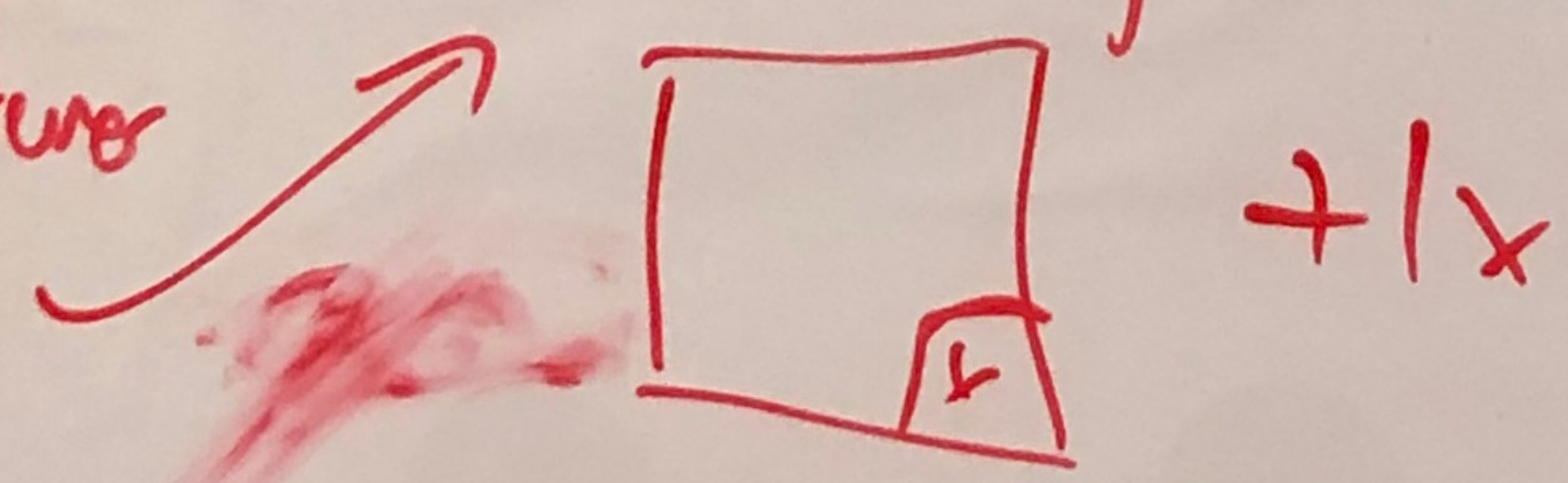
W D ≠ I B a w @ s e n t @ n s ?

S =

1 0

Sentence length

in whatever units



prioritise rarer diphones

Max diphones not in db.
Context
prosodic position
sentence position.
Word position
Stress.

Step 1: Sources of Text

- NEWS Articles (Domain: Sports)
 - BBC ↪ and transcripts (radio news)
 - AP
 - NPR
- TED Talks transcripts

Step 2: Clean the Text

- Split punctuation
- remove sentences that are too long, duplicate words
- tokenize
- normalize (eg. Dr. , 1-0 score)

Step 3: Richness Measure

this is a laptop
dhis iz @ laeptap

- new diphones normalised by ^{no. of} diphones
 - reverse Zipf distribution
- diphones in new context
 - position in word
 - position in utterance
- within question: preceding 3 diphones or whole sentence?

STEP 1

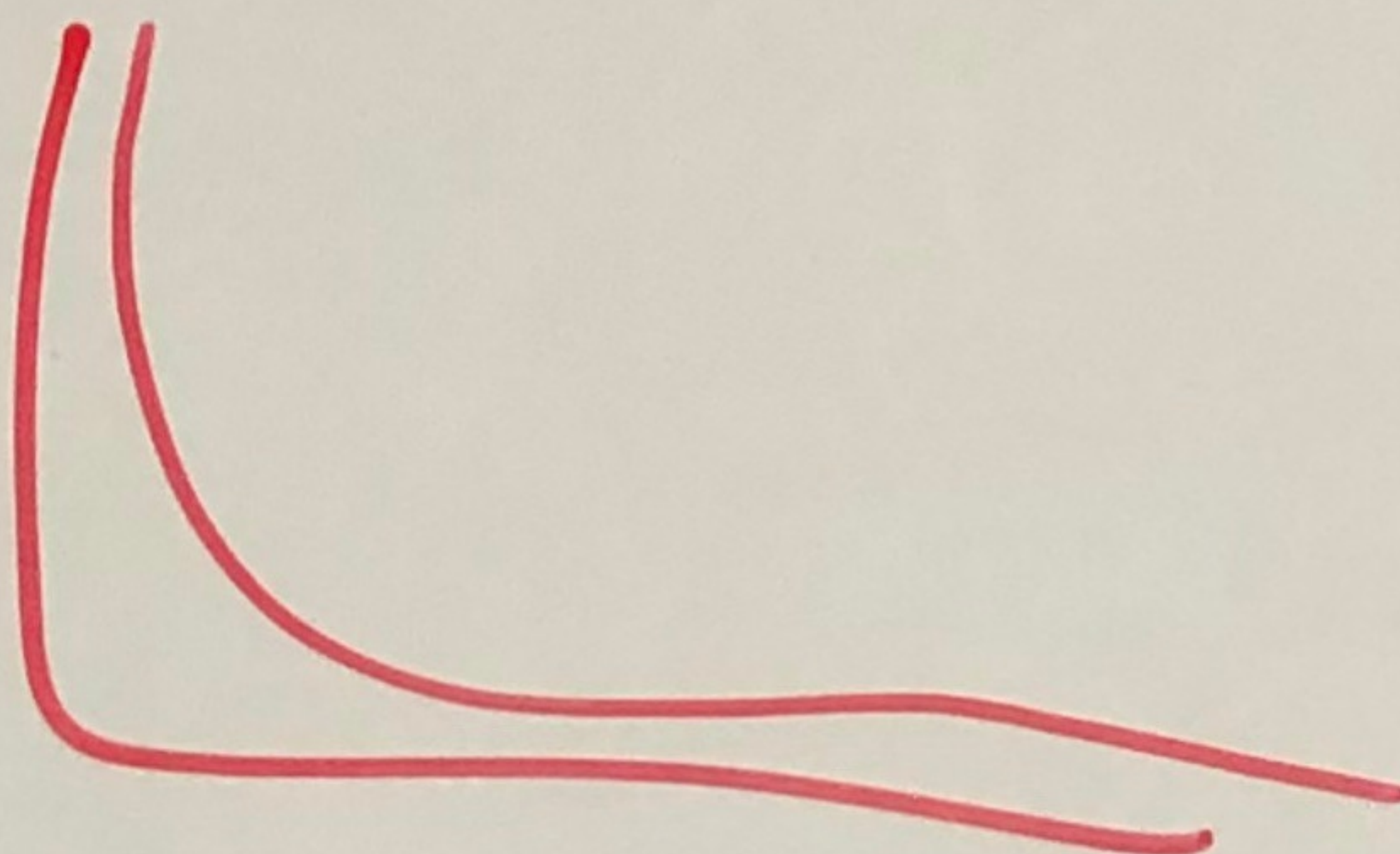
SOURCE: articles, wikipedia, GPT, podcast.
transcripts, film transcripts?
↳ theatre

STEP 2

CLEAN: pre-process text (remove NSWs?, segment by punctuation, discard every sentence w/ symbols, remove sentences w/ words not in dictionary, keep sentences of length 5-15 words,

~~the~~ cat slept on ~~the~~

[the



]

domain specific vocab
embeddings

①

radio shows / podcasts
news transcriptions copyright?

british national corpus

CNN, daily mail

~~scribble~~ Scribe (british national corpus)

weatherforecast corpus

②

beautiful soup
for not-friendly
formatted text

~~scribble~~ FLEXH
- grab
the readability

(hat gpt)

use Language model
to calculate perplexity

nlTK - tokenize, lower case
~~scribble~~ split into
sentences

→ NSWs (ver...)

transcribed things were
designed to be spoken,
should be easy to read

③

weighted average of context

questions (prosody)
stress
phrase initial/final

run text through G2P
to find all diphones for
maximal coverage

BPE
to find
common
morphemes

prosody
↳ through
punctuation
! ?

prosody + diphone coverage