# Linear Text Segmentation with Neural OIE on Novels and Subtitles

Olena Mashkina

# Do Androids Dream of Text Segmentation?

# Outline

1. Related work
2. Dataset
3. Research questions
4. Method
5. Results

# Related work

# Linear text segmentation

| algorithm | publication year | supervised | similarity-based | generative | key features |
|---|---|---|---|---|---|
| TextTiling | 1997 | no | ✓ | | lexical co-occurance |
| C99 | 2000 | no | ✓ | | ranking matrix, divisive hierarchical clustering |
| U00 | 2001 | no | | ✓ | minimum cost segmentation, dynamic programming |
| LCSeg | 2003 | no | ✓ | | TextTiling-based, lexical chains |
| Sun et al. | 2007 | no | | ✓ | mutual information, dynamic programming |
| BayesSeg | 2008 | no | | ✓ | Bayesian framework, incorporating cue phrases, dynamic programming |
| TopicTiling | 2012 | yes | | ✓ | LDA-based |
| GraphSeg | 2016 | no | ✓ | | semantic relatedness graph, word embeddings |
| Sector | 2019 | no | | | LSTM neural network, topic labeling |

# Open Information Extraction

| OIE system | publication year | key features |
|---|---|---|
| TextRunner | 2007 | <ul><li>POS tags, NP chunks</li><li>labeling by handcrafted patterns</li></ul> |
| ReVerb | 2011 | <ul><li>rule-based</li><li>POS tags, NP chunks</li><li>lexical and semantic constraints</li></ul> |
| ClausIE | 2013 | <ul><li>clause-based</li><li>dependency parsing</li><li>no training required</li></ul> |
| Stanford OIE | 2015 | <ul><li>clause-based</li><li>dependency parsing</li><li>minimization of extracted clauses</li><li>extraction based on handcrafted patterns</li></ul> |
| RnnOIE | 2018 | <ul><li>neural-based</li><li>bi-LSTM transducer for supervised model training</li><li>extracts n-ary relational tuples</li><li>OIE as a sequence labeling problem</li></ul> |
| Cui et al. | 2018 | <ul><li>neural-based</li><li>encoder-decoder LSTM RNN for supervised model training</li><li>extracts binary relational tuples</li><li>OIE as a sequence-to-sequence generation problem</li></ul> |

# Dataset

# Dataset (novels)

Selection criteria:

- dystopian genre
- structured in chapters
- adapted into film
- in English

|    | Novel title | Author | Publication year | Length in words |
|----|-------------|--------|------------------|-----------------|
| 1  | 1984: A Novel | G. Orwell | 1949 | 101,327 |
| 2  | Brave New World | A. Huxley | 1932 | 66,511 |
| 3  | We | Y. Zamyatin | Written in 1920, translated into English in 1924 | 62,794 |
| 4  | The Handmaid's Tale | M. Atwood | 1985 | 94,643 |
| 5  | Do Androids Dream of Electric Sheep? | P. K. Dick | 1968 | 64,106 |
| 6  | The Hunger Games | S. Collins | 2008 | 103,624 |
| 7  | Catching Fire | S. Collins | 2009 | 105,631 |
| 8  | Mockingjay | S. Collins | 2010 | 104,812 |
| 9  | The Giver | L. Lowry | 1993 | 44,790 |
| 10 | The Maze Runner | J. Dashner | 2009 | 79,431 |
| 11 | Ready Player One | E. Cline | 2011 | 140,721 |

# Dataset (subtitles)

Selection criteria:

- film is an adaptation of a dystopian novel
- subtitle-novel pair
- in English

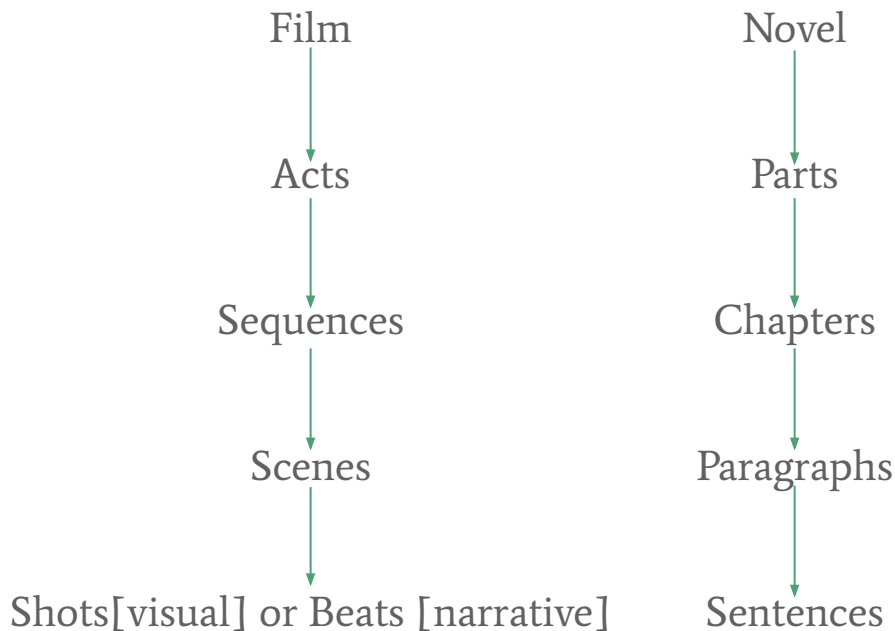| | Film title | Director | Release year | Length in words |
|---|---|---|---|---|
| 1 | 1984 | M. Anderson | 1956 | 6,252 |
| 2 | 1984 | M. Radford | 1984 | 7,011 |
| 3 | Brave New World | L. Libman, L. Williams | 1998 | 7,464 |
| 4 | We | V. Jasný | 1982 | 6,587 |
| 5 | The Handmaid's Tale | V. Schlöndorff | 1990 | 6,448 |
| 6 | Blade Runner | R. Scott | 1982 | 4,303 |
| 7 | The Hunger Games | G. Ross | 2012 | 6,365 |
| 8 | The Hunger Games: Catching Fire | F. Lawrence | 2013 | 9,309 |
| 9 | The Hunger Games: Mockingjay - Part 1 | F. Lawrence | 2014 | 8,443 |
| 10 | The Hunger Games: Mockingjay - Part 2 | F. Lawrence | 2015 | 8,517 |
| 11 | The Giver | P. Noyce | 2014 | 7,172 |
| 12 | The Maze Runner | W. Ball | 2014 | 6,523 |
| 13 | Ready Player One | S. Spielberg | 2018 | 10,863 |

# Research questions

# Research question

**Related to Natural Language Processing (NLP):**

How does a modification of a linear text segmentation method by adding word embeddings and knowledge generated by Open Information Extraction (OIE) influence the performance of this method?

# Research question

**Related to the created dataset:**

How does the performance
of presented in this work text
segmentation pipeline compare
for different fictional narrative text
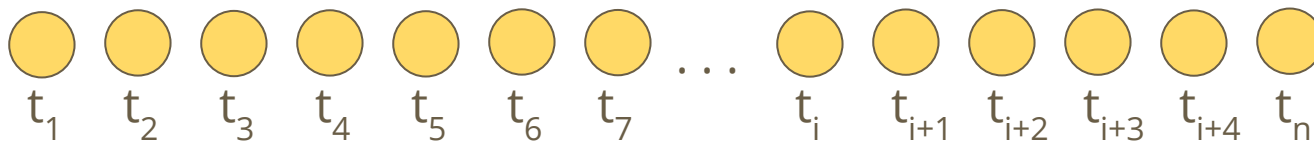corpora (novels and subtitles)?

Film
↓
Acts
↓
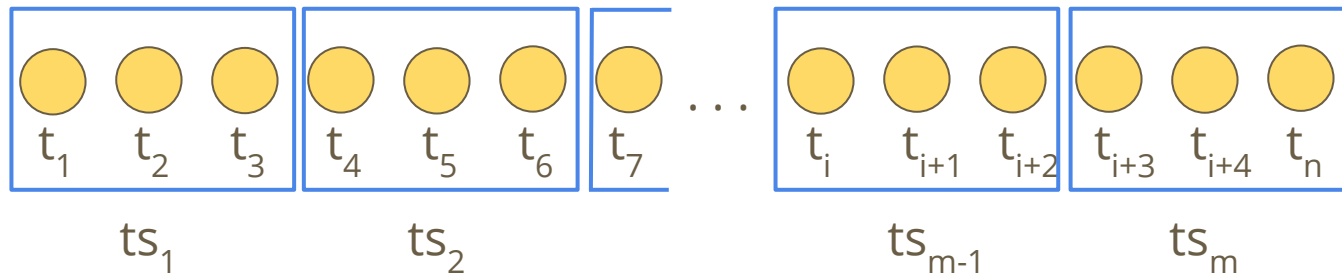Sequences
↓
Scenes
↓
Shots[visual] or Beats [narrative]

Novel
↓
Parts
↓
Chapters
↓
Paragraphs
↓
Sentences

# Method

# TextTiling

# TextTiling

- words of the input text are lemmatized
- a series of word tokens $t_1$ ... $t_n$



$t_1$  $t_2$  $t_3$  $t_4$  $t_5$  $t_6$  $t_7$  ...  $t_i$  $t_{i+1}$  $t_{i+2}$  $t_{i+3}$  $t_{i+4}$  $t_n$

M. A. Hearst. Texttiling: Segmenting text into multi-paragraph subtopic passages. Journal of Computational linguistics, 23(1):33–64, 1997

# TextTiling

- token-sequence size *w*
- w approximates the length of the sentence
- *w* = 3 in this example

$$ts_1 \quad [t_1 \; t_2 \; t_3] \qquad ts_2 \quad [t_4 \; t_5 \; t_6] \qquad [t_7] \; \ldots \quad ts_{m-1} \quad [t_i \; t_{i+1} \; t_{i+2}] \qquad ts_m \quad [t_{i+3} \; t_{i+4} \; t_n]$$

# TextTiling

- token-sequence size $w$
- $w$ approximates the length of the sentence
- $w = 3$ in this example



$ts_1$    $ts_2$    $ts_{m-1}$    $ts_m$
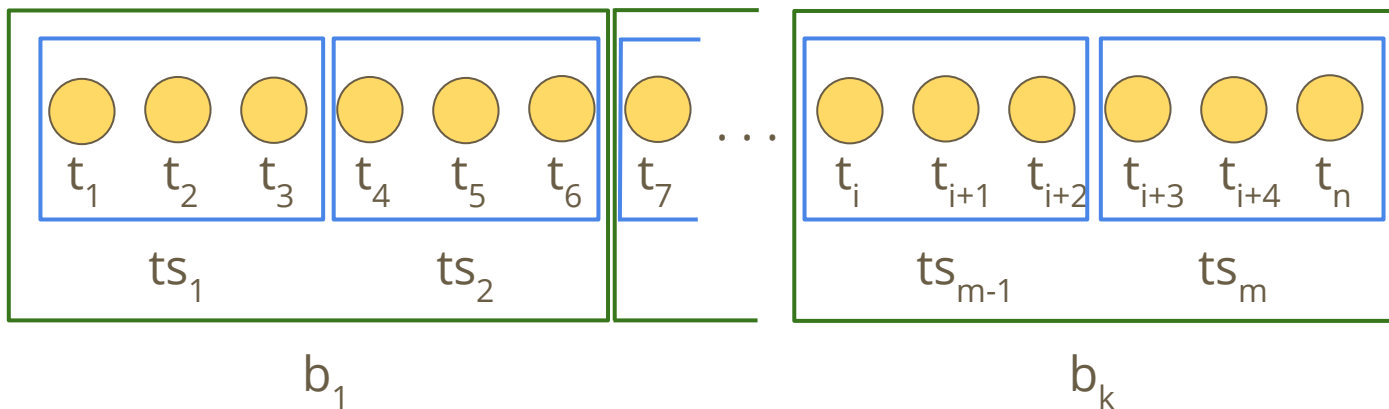
token-sequence gap

# TextTiling

- block size k
- block size approximates the length of the paragraph in sentences
- k = 2 in this example

# TextTiling (lexical score)

- vocabulary change signifies a change of subtopic in text
- a lexical score is computed between 2 neighboring blocks at each step
- moving window: shift by one token-sequence, compare 2 blocks
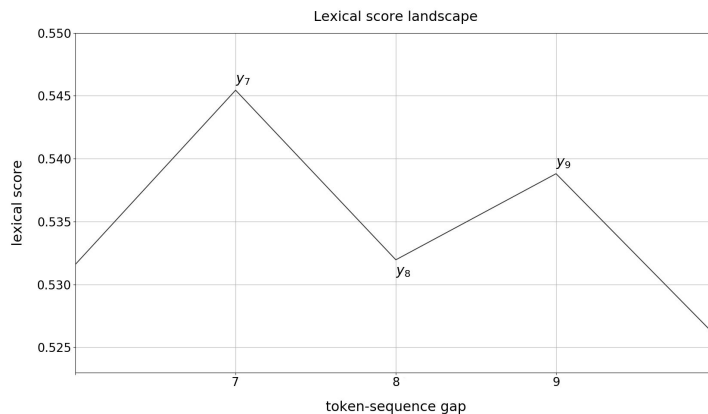- each token-sequence gap is assigned a lexical score

# TextTiling (lexical score)

- lexical score value is cosine similarity between blocks

$$score = \frac{\sum_t w_{t,b1} w_{t,b2}}{\sqrt{\sum_t w_{t,b1}^2 w_{t,b2}^2}}$$

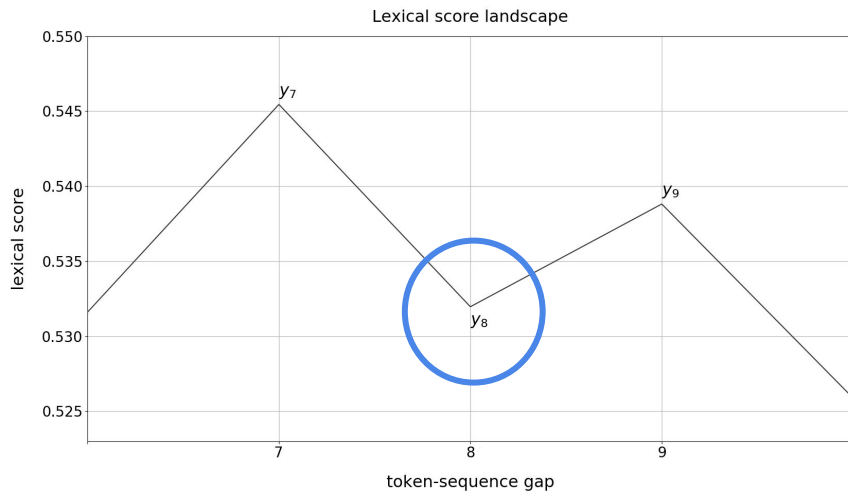$w_{t,b1}$ ... frequency of a vocabulary token $t$ within a block $b_1$

$w_{t,b2}$ ... frequency of a vocabulary token $t$ within a block $b_2$



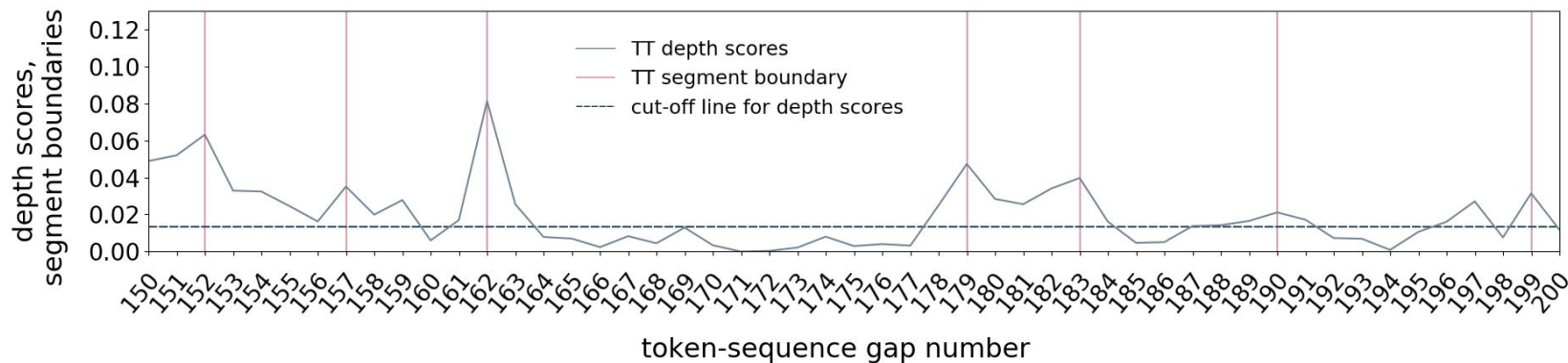Lexical score landscape

# TextTiling (depth score)

- the depth of the valley

- distance from the top peaks on both sides of the token-sequence gap

$$depth\_score(gap_j) = score(gap_l) - score(gap_j) + score(gap_r) - score(gap_j)$$

# TextTiling (depth score)

- the larger the depth score value, the more probable is a topic switch

# Word embeddings

# Word embeddings (model)

- a novel and its film adaptation share the fictional terms and proper names

- a model was trained on the novel and fine-tuned on subtitle text

- stop words were not included as they do not carry semantic meaning and should have no significant impact on the vector space

- the same model was used at the parameter optimization step

# Word embeddings

- at the step of splitting input into token sequences word tokens are replaced by their vector representations

- a block is represented by a vector (unchanged)

- moving window (unchanged): shift by one token-sequence, calculate lexical score of neighboring blocks

- a block is represented by a vector sum of all word embedding vectors in the block

# Open Information Extraction

# Open Information Extraction weights

- propositions are extracted from the sentence in form of n-ary relational tuples

- a word token may not be a part of any tuple and multiple tuples may share the same word token

- frequency of a word token in tuples corresponds to the strength of syntactic meaning

- an overall number of occurrences of a word token in all extracted propositions of a single sentence is considered token's weight

# Open Information Extraction weights (example)

Sentence: Better her than me, Rita said, and I opened the door.

Proposition 1: "Better her than me , [ARG0: Rita] [V: said] , and I opened the door ."

Proposition 2: "Better her than me , Rita said , and [ARG0: I] [V: opened] [ARG1: the door] ."

| | Better | her | than | me | Rita | said | and | I | opened | the | door |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Proposition 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| Proposition 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| Term weight | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 1 |

# Results

# Results (novels)

- replacing word tokens by their embedding vector representations

- decrease in the WindowDiff measure

- 9 out of 11 novels

| | Novel title | Parameters | TextTiling | **TextTiling with word embeddings** |
|---|---|---|---|---|
| 1 | 1984: A Novel | k=24, w=136 | WD=0.961 $P_k$=0.531 | **WD=0.917** $P_k$=0.514 |
| 2 | Brave New World | k=100, w=258 | WD=0.542 $P_k$=0.424 | WD=0.623 $P_k$=0.503 |
| 3 | We | k=22, w=88 | WD=0.783 $P_k$=0.521 | **WD=0.703** $P_k$=0.499 |
| 4 | The Handmaid's Tale | k=14, w=112 | WD=0.788 $P_k$=0.525 | **WD=0.721** $P_k$=0.485 |
| 5 | Do Androids Dream of Electric Sheep? | k=20, w=70 | WD=0.99 $P_k$=0.523 | **WD=0.953** $P_k$=0.508 |
| 6 | The Hunger Games | k=39, w=54 | WD=0.996 $P_k$=0.509 | **WD=0.981** $P_k$=0.512 |
| 7 | Catching Fire | k=11, w=44 | WD=1.0 $P_k$=0.509 | WD=1.0 $P_k$=0.509 |
| 8 | Mockingjay | k=19, w=46 | WD=1.0 $P_k$=0.524 | **WD=0.989** $P_k$=0.528 |
| 9 | The Giver | k=12, w=107 | WD=0.723 $P_k$=0.529 | **WD=0.69** **$P_k$=0.468** |
| 10 | The Maze Runner | k=13, w=49 | WD=0.989 $P_k$=0.511 | **WD=0.947** $P_k$=0.503 |
| 11 | Ready Player One | k=8, w=71 | WD=0.998 $P_k$=0.51 | **WD=0.992** $P_k$=0.513 |

# Results (subtitles)

- replacing word tokens by their embedding vector representations

- decrease in the WindowDiff measure

- 6 out of 13 subtitle files

| | Film title | Parameters | TextTiling | TextTiling with word embeddings |
|---|---|---|---|---|
| 1 | 1984 (1956) | k=39, w=22 | WD=0.412 $P_k$=0.387 | WD=0.423 $P_k$=0.401 |
| 2 | 1984 (1984) | k=29, w=22 | WD=0.41 $P_k$=0.393 | **WD=0.395** $P_k$=0.376 |
| 3 | Brave New World | k=32, w=41 | WD=0.395 $P_k$=0.372 | **WD=0.393** $P_k$=0.37 |
| 4 | We | k=50, w=48 | WD=0.393 $P_k$=0.385 | WD=0.413 $P_k$=0.408 |
| 5 | The Handmaid's Tale | k=56, w=44 | WD=0.349 $P_k$=0.334 | WD=0.367 $P_k$=0.352 |
| 6 | Blade Runner | k=18, w=34 | WD=0.336 $P_k$=0.325 | WD=0.36 $P_k$=0.351 |
| 7 | The Hunger Games | k=46, w=20 | WD=0.374 $P_k$=0.359 | WD=0.411 $P_k$=0.4 |
| 8 | The Hunger Games: Catching Fire | k=57, w=40 | WD=0.32 $P_k$=0.305 | WD=0.334 $P_k$=0.324 |
| 9 | The Hunger Games: Mockingjay - Part 1 | k=13, w=97 | WD=0.354 $P_k$=0.35 | **WD=0.335** $P_k$=0.331 |
| 10 | The Hunger Games: Mockingjay - Part 2 | k=21, w=97 | WD=0.326 $P_k$=0.322 | **WD=0.315** $P_k$=0.312 |
| 11 | The Giver | k=,82 w=28 | WD=0.427 $P_k$=0.388 | **WD=0.408** $P_k$=0.372 |
| 12 | The Maze Runner | k=43, w=27 | WD=0.357 $P_k$=0.338 | WD=0.373 $P_k$=0.36 |
| 13 | Ready Player One | k=95, w=34 | WD=0.365 $P_k$=0.346 | **WD=0.354** $P_k$=0.339 |

# Results (novels)

- replacing word tokens by their embedding vector representations

- OIE weights

- decrease in the WindowDiff measure

- 7 out of 11 novels

| | Novel title | Parameters | TextTiling | TextTiling with word embeddings and OIE weights |
|---|---|---|---|---|
| 1 | 1984: A Novel | k=24, w=136 | WD=0.961 $P_k$=0.531 | **WD=0.854** $P_k$=0.515 |
| 2 | Brave New World | k=100, w=258 | WD=0.542 $P_k$=0.424 | WD=0.685 $P_k$=0.541 |
| 3 | We | k=22, w=88 | WD=0.783 $P_k$=0.521 | **WD=0.739** $P_k$=0.5 |
| 4 | The Handmaid's Tale | k=14, w=112 | WD=0.788 $P_k$=0.525 | **WD=0.709** $P_k$=0.518 |
| 5 | Do Androids Dream of Electric Sheep? | k=20, w=70 | WD=0.99 $P_k$=0.523 | **WD=0.983** $P_k$=0.517 |
| 6 | The Hunger Games | k=39, w=54 | WD=0.996 $P_k$=0.509 | **WD=0.9814** $P_k$=0.515 |
| 7 | Catching Fire | k=11, w=44 | WD=1.0 $P_k$=0.509 | WD=1.0 $P_k$=0.509 |
| 8 | Mockingjay | k=19, w=46 | WD=1.0 $P_k$=0.524 | WD=1.0 $P_k$=0.524 |
| 9 | The Giver | k=12, w=107 | WD=0.723 $P_k$=0.529 | WD=0.739 $P_k$=0.481 |
| 10 | The Maze Runner | k=13, w=49 | WD=0.989 $P_k$=0.511 | **WD=0.934** **$P_k$=0.498** |
| 11 | Ready Player One | k=8, w=71 | WD=0.998 $P_k$=0.51 | **WD=0.981** $P_k$=0.516 |

# Results (subtitles)

- replacing word tokens by their embedding vector representations

- OIE weights

- decrease in the WindowDiff measure

- 4 out of 13 subtitle files

| | Film title | Parameters | TextTiling | TextTiling with word embeddings and OIE weights |
|---|---|---|---|---|
| 1 | 1984 (1956) | k=39, w=22 | WD=0.412 $P_k$=0.387 | **WD=0.389** $P_k$=0.375 |
| 2 | 1984 (1984) | k=29, w=22 | WD=0.41 $P_k$=0.393 | **WD=0.405** $P_k$=0.389 |
| 3 | Brave New World | k=32, w=41 | WD=0.395 $P_k$=0.372 | **WD=0.381** $P_k$=0.357 |
| 4 | We | k=50, w=48 | WD=0.393 $P_k$=0.385 | WD=0.445 $P_k$=0.445 |
| 5 | The Handmaid's Tale | k=56, w=44 | WD=0.349 $P_k$=0.334 | WD=0.369 $P_k$=0.36 |
| 6 | Blade Runner | k=18, w=34 | WD=0.336 $P_k$=0.325 | WD=0.354 $P_k$=0.343 |
| 7 | The Hunger Games | k=46, w=20 | WD=0.374 $P_k$=0.359 | WD=0.411 $P_k$=0.394 |
| 8 | The Hunger Games: Catching Fire | k=57, w=40 | WD=0.32 $P_k$=0.305 | WD=0.324 $P_k$=0.327 |
| 9 | The Hunger Games: Mockingjay - Part 1 | k=13, w=97 | WD=0.354 $P_k$=0.35 | **WD=0.349** $P_k$=0.347 |
| 10 | The Hunger Games: Mockingjay - Part 2 | k=21, w=97 | WD=0.326 $P_k$=0.322 | WD=0.329 $P_k$=0.326 |
| 11 | The Giver | k=,82 w=28 | WD=0.427 $P_k$=0.388 | WD=0.442 $P_k$=0.414 |
| 12 | The Maze Runner | k=43, w=27 | WD=0.357 $P_k$=0.338 | WD=0.38 $P_k$=0.367 |
| 13 | Ready Player One | k=95, w=34 | WD=0.365 $P_k$=0.346 | WD=0.378 $P_k$=0.369 |

# Findings (general)

- generalization of TextTiling parameters which would satisfy all input files equally well is not possible

- automatically generated ground truth may be too coarse for novels

- dialogues present a big challenge for the pipeline

# Findings (research questions)

- pipeline is more effective for novels than subtitles (grammatically incomplete informal sentences)

- word embedding have a potential to improve the performance of TextTiling (increase in performance for 6 out of 13 subtitles and 9 out of 11 novels)

- application of word embeddings and OIE weights has a potential to improve the performance of TextTiling (increase in performance for 4 out of 13 subtitles and 7 out of 11 novels)

# Thank you for your attention

# Additional slides

# Basic concepts

# Basic concepts: linear text segmentation

Goal: to automatically locate a transition from one topic to another in a text

Result:

- text is separated into non-overlapping neighboring textual segments
- each segment characterized by a single homogeneous topic
- each segment contains a certain number of passages (e.g. paragraphs or sentences)

# Basic concepts: linear text segmentation example

It was, he now realised, because of this other incident that he had suddenly decided to come home and begin with the diary today.

text segment 1

It had happened that morning at the Ministry, if anything so nebulous could be said to happen.

text segment 2

Source: Nineteen Eighty-Four by G. Orwell

# Basic concepts: Open Information Extraction

Goal: to create a representation of propositions in a text document in form of n-tuples

Result:

- each sentence in the document is assigned a set of relational n-tuples
- each tuple contains at least two arguments connected by a semantic relation between them (predicate): {argument 1, predicate, argument 2}
- n-tuples should represent propositions clearly expressed in the sentence
- there is no limit to the number of tuples extracted from a single sentence

# Basic concepts: Open Information Extraction example

Input sentence:

The sun went down and the dark-gray clouds changed color.

Extracted propositions:

1) [ARG0: The sun] [V: went] [ARG1: down]

2) [ARG0: the dark - gray clouds] [V: changed] [ARG1: color]

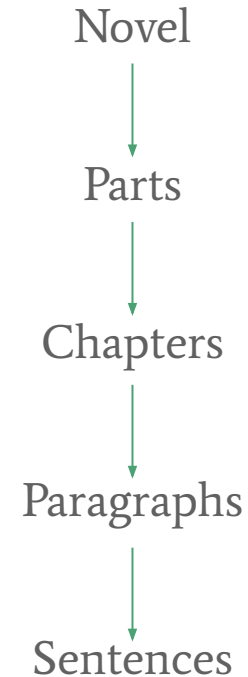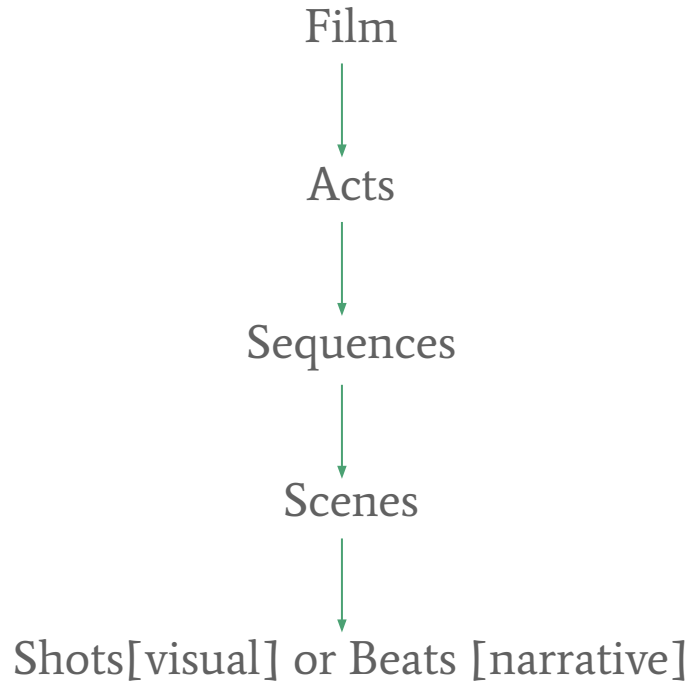# Basic concepts: word embeddings

Goal:

transforming original textual data into a vector space based on prediction from the linguistic context
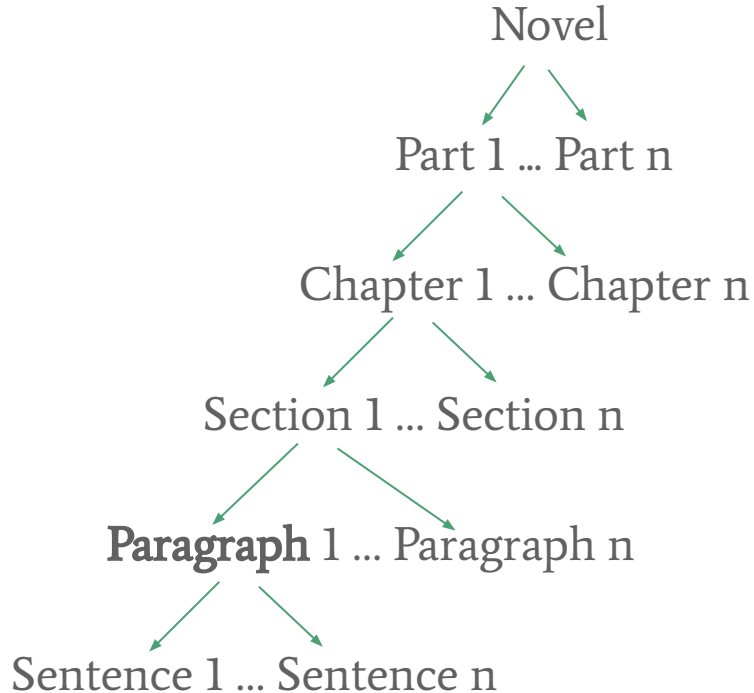
Result:

- each word in a vocabulary of a corpus is assigned a single real-valued vector
- approximation of word's meaning. Distributional hypothesis: words occurring in similar context tend to have similar meaning
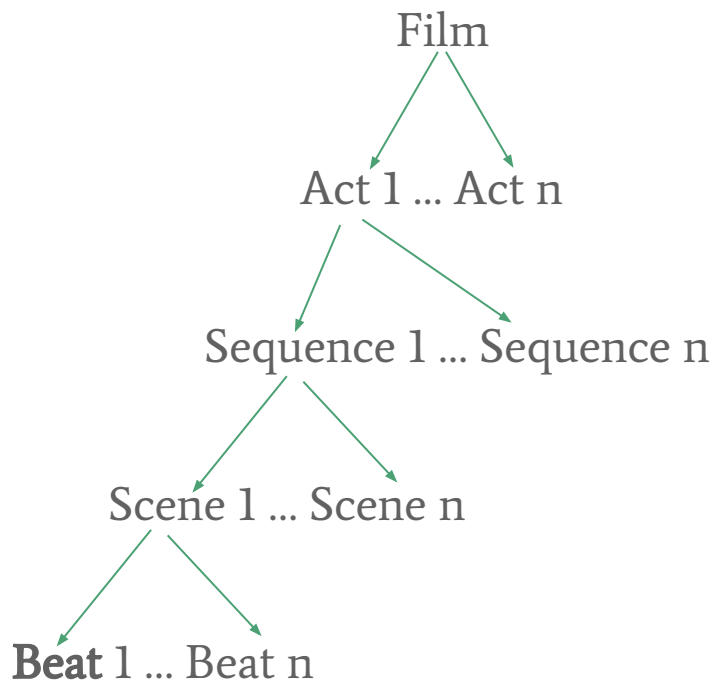
# Narrative structure

# Narrative structure

Film

↓

Acts

↓

Sequences

↓

Scenes

↓

Shots[visual] or Beats [narrative]

Novel

↓

Parts

↓

Chapters

↓

Paragraphs

↓

Sentences

# Novel structure

Novel

↙ ↘

Part 1 ... Part n

↙ ↘

Chapter 1 ... Chapter n

↙ ↘

Section 1 ... Section n

↙ ↘

**Paragraph** 1 ... Paragraph n

↙ ↘

Sentence 1 ... Sentence n

A paragraph is

- a subdivision of a written composition
- begins on a new usually indented line
- consists of one or more sentences
- deals with one point or
- gives the words of one speaker

# Film screenplay structure

Film

Act 1 ... Act n

Sequence 1 ... Sequence n

Scene 1 ... Scene n

**Beat** 1 ... Beat n

A beat

- is an action/reaction event for moving the plot forward
- should stimulate an emotion from the audience

https://tvtropes.org/pmwiki/pmwiki.php/Main/NarrativeBeats

# Dataset

# Dataset: ground truth for novels

Chapter 1

A squat grey building of only thirty-four stories.

   .        .        .

   .        .        .

   .        .        .

"Just one glance."

Chapter 2

Mr. Foster was left in the Decanting Room.

   .        .        .

   .        .        .

   .        .        .

# Dataset: ground truth for novels

~~Chapter 1~~

A squat grey building of only thirty-four stories.

   .         .         .
   .         .         .
   .         .         .

"Just one glance."

one text segment

~~Chapter 2~~

Mr. Foster was left in the Decanting Room.

   .         .         .
   .         .         .
   .         .         .

- each chapter is treated as a single text segment
- headers are filtered out

# Dataset: ground truth for subtitles

77
01:28:03,196 --> 01:28:06,108
Ok. Checked and cleared.
Have a better one.

11.35 seconds

778
01:28:17,460 --> 01:28:19,262
Hello.
Hi. Is J.F. there?

779
01:28:19,263 --> 01:28:20,284
Who is it?

780
01:28:20,485 --> 01:28:22,149
This is Eddie, old friend of J.F.'s.

~4.78 seconds

781
01:28:26,928 --> 01:28:28,638
That's no way to treat a friend.

782
01:28:31:27,233 --> 01:31:29,318
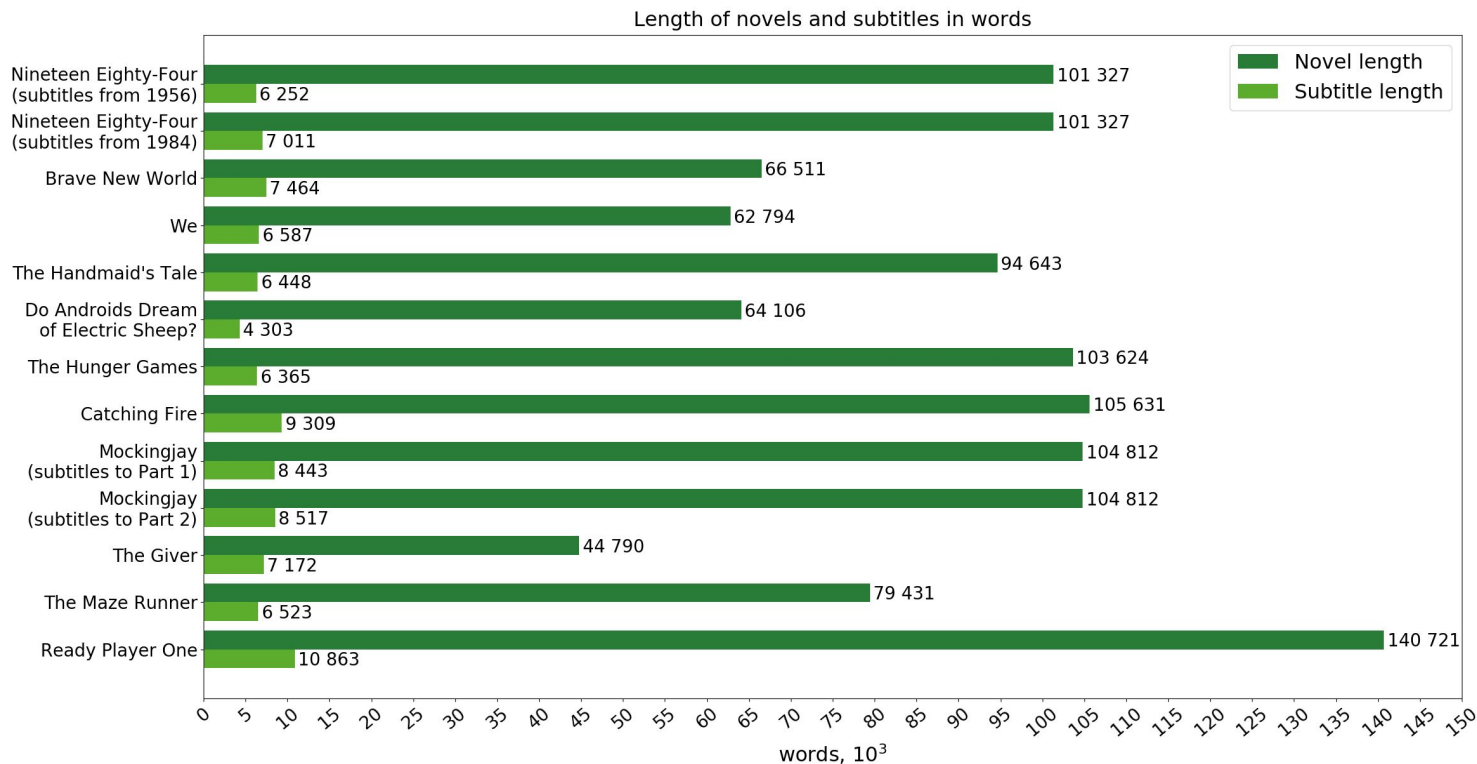Home again, home again, jiggidy-jig.

178.595 seconds

The beginning of a new text segment is identified if there was a pause between subtitle sequences longer than 5 seconds.

one text segment:

"Hello. Hi. Is J.F. there? Who is it? This is Eddie, old friend of J.F.'s. That's no way to treat a friend."

Source: Blade Runner (1982)

# Dataset motivation



Length of novels and subtitles in words

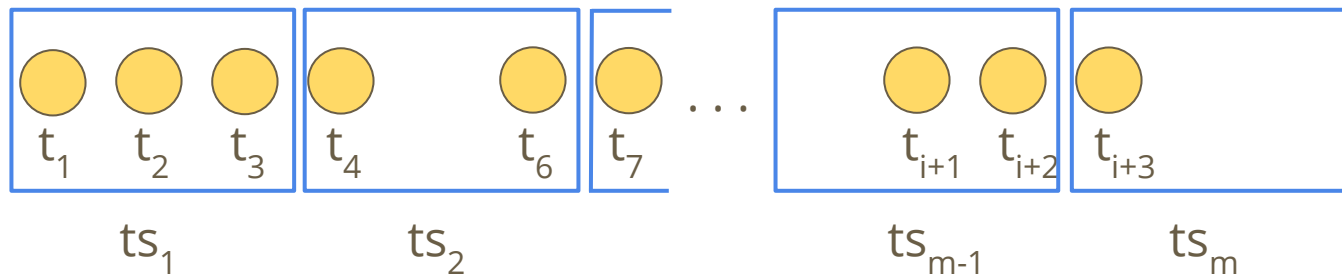| Book | Novel length | Subtitle length |
|---|---|---|
| Nineteen Eighty-Four (subtitles from 1956) | 101 327 | 6 252 |
| Nineteen Eighty-Four (subtitles from 1984) | 101 327 | 7 011 |
| Brave New World | 66 511 | 7 464 |
| We | 62 794 | 6 587 |
| The Handmaid's Tale | 94 643 | 6 448 |
| Do Androids Dream of Electric Sheep? | 64 106 | 4 303 |
| The Hunger Games | 103 624 | 6 365 |
| Catching Fire | 105 631 | 9 309 |
| Mockingjay (subtitles to Part 1) | 104 812 | 8 443 |
| Mockingjay (subtitles to Part 2) | 104 812 | 8 517 |
| The Giver | 44 790 | 7 172 |
| The Maze Runner | 79 431 | 6 523 |
| Ready Player One | 140 721 | 10 863 |

words, $10^3$

# Benefits of the dataset

- synonyms rather than word repetitions
- made up terms
- already existing words may obtain ironic meaning (e.g. "Ministry of Love")
- irregular length of sentences and paragraphs
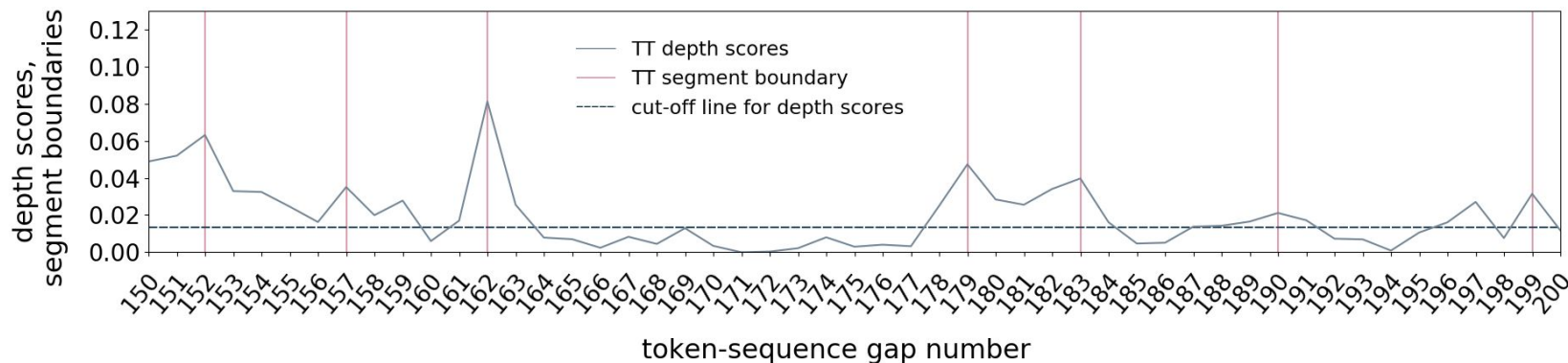- unconfined text structure

# Method

# TextTiling

- stop words are removed

# TextTiling (depth score)

- depth scores are sorted (the highest depth score is a guaranteed boundary)

- we want to avoid many boundaries very close to each other so at least three token sequences are required between boundaries
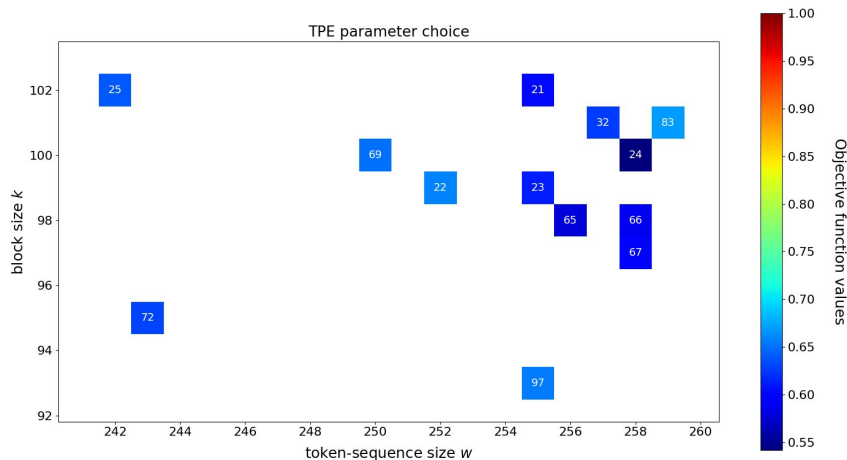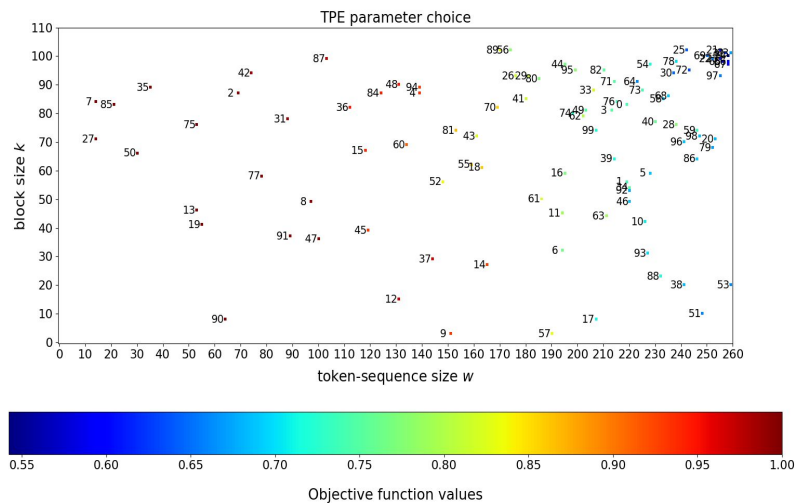
# Parameter search

- the minimum of the token-sequence value was set to median sentence length

- the maximum of the token-sequence value was set to maximum sentence length

- the minimum of the block size value was set to median paragraph length in sentences

- the maximum of the block size value was set to maximum paragraph length in sentences

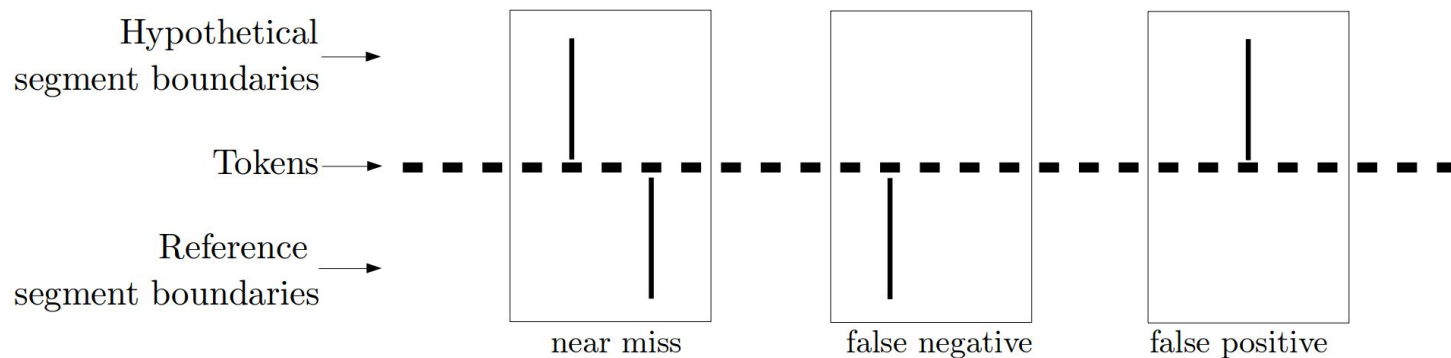- the output of objective function to minimize is WindowDiff measure value

# Parameter search

$$\forall w \in [median\_sl, max\_sl], length\ in\ tokens$$

$$\forall k \in [median\_pl, max\_pl], length\ in\ sentences$$

# Evaluation metrics

$P_k$ does not penalize if the number of hypothetical boundaries exceeds the number of reference boundaries in the window

# Results

# Results (novels)

- word embeddings:

  increased performance for 9 out of 11 novels compared to TextTiling

- word embeddings with OIE weights:

  increased performance for 7 out of 11 novels compared to TextTiling

| | Novel title | Parameters | TextTiling | TextTiling with word embeddings | TextTiling with OIE weights |
|---|---|---|---|---|---|
| 1 | 1984: A Novel | k=24, w=136 | WD=0.961 $P_k$=0.531 | WD=0.917 $P_k$=0.514 | **WD=0.854** $P_k$=0.515 |
| 2 | Brave New World | k=100, w=258 | **WD=0.542** $P_k$=0.424 | WD=0.623 $P_k$=0.503 | WD=0.685 $P_k$=0.541 |
| 3 | We | k=22, w=88 | WD=0.783 $P_k$=0.521 | **WD=0.703** $P_k$=0.499 | WD=0.739 $P_k$=0.5 |
| 4 | The Handmaid's Tale | k=14, w=112 | WD=0.788 $P_k$=0.525 | WD=0.721 $P_k$=0.485 | **WD=0.709** $P_k$=0.518 |
| 5 | Do Androids Dream of Electric Sheep? | k=20, w=70 | WD=0.99 $P_k$=0.523 | **WD=0.953** $P_k$=0.508 | WD=0.983 $P_k$=0.517 |
| 6 | The Hunger Games | k=39, w=54 | WD=0.996 $P_k$=0.509 | **WD=0.981** $P_k$=0.512 | WD=0.9814 $P_k$=0.515 |
| 7 | Catching Fire | k=11, w=44 | WD=1.0 $P_k$=0.509 | WD=1.0 $P_k$=0.509 | WD=1.0 $P_k$=0.509 |
| 8 | Mockingjay | k=19, w=46 | WD=1.0 $P_k$=0.524 | **WD=0.989** $P_k$=0.528 | WD=1.0 $P_k$=0.524 |
| 9 | The Giver | k=12, w=107 | WD=0.723 $P_k$=0.529 | **WD=0.69** **$P_k$=0.468** | WD=0.739 $P_k$=0.481 |
| 10 | The Maze Runner | k=13, w=49 | WD=0.989 $P_k$=0.511 | WD=0.947 $P_k$=0.503 | **WD=0.934** **$P_k$=0.498** |
| 11 | Ready Player One | k=8, w=71 | WD=0.998 $P_k$=0.51 | WD=0.992 $P_k$=0.513 | **WD=0.981** $P_k$=0.516 |

# Results (subtitles)

- word embeddings:

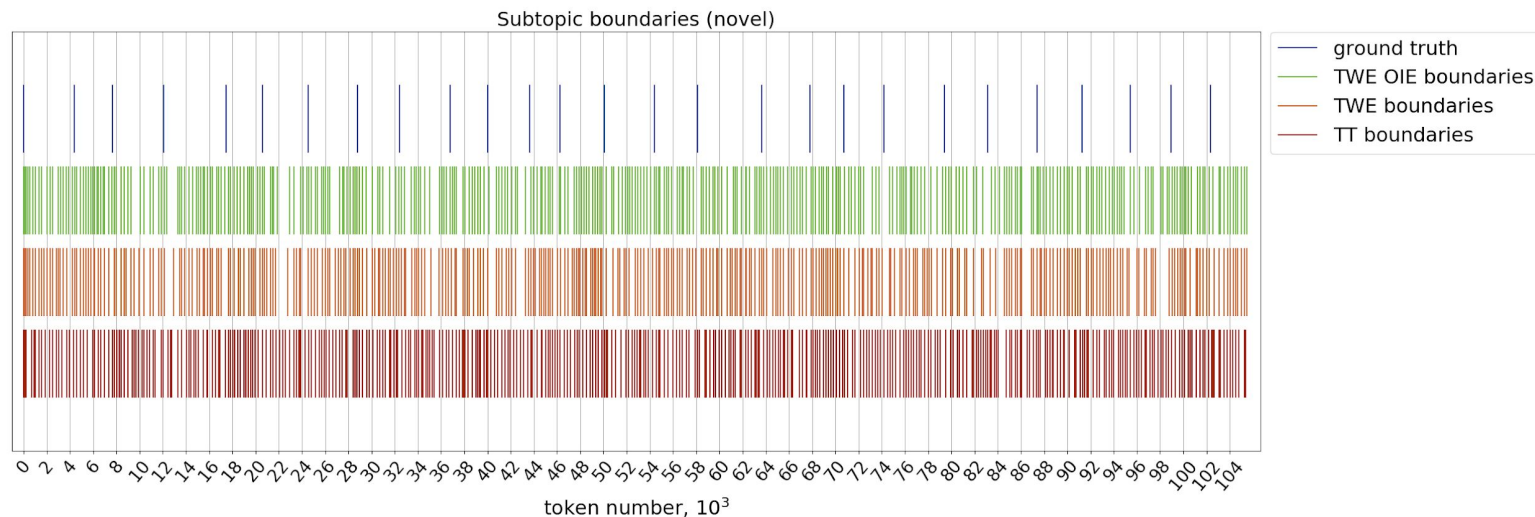  increased performance for 6 out of 13 subtitle files compared to TextTiling

- word embeddings with OIE weights:

  increased performance for 4 out of 13 subtitle files compared to TextTiling

|  | Film title | Parameters | TextTiling | TextTiling with word embeddings | TextTiling with OIE weights |
|---|---|---|---|---|---|
| 1 | 1984 (1956) | k=39, w=22 | WD=0.412 $P_k$=0.387 | WD=0.423 $P_k$=0.401 | **WD=0.389** $P_k$=0.375 |
| 2 | 1984 (1984) | k=29, w=22 | WD=0.41 $P_k$=0.393 | **WD=0.395** $P_k$=0.376 | WD=0.405 $P_k$=0.389 |
| 3 | Brave New World | k=32, w=41 | WD=0.395 $P_k$=0.372 | WD=0.393 $P_k$=0.37 | **WD=0.381** $P_k$=0.357 |
| 4 | We | k=50, w=48 | **WD=0.393** $P_k$=0.385 | WD=0.413 $P_k$=0.408 | WD=0.445 $P_k$=0.445 |
| 5 | The Handmaid's Tale | k=56, w=44 | **WD=0.349** $P_k$=0.334 | WD=0.367 $P_k$=0.352 | WD=0.369 $P_k$=0.36 |
| 6 | Blade Runner | k=18, w=34 | **WD=0.336** $P_k$=0.325 | WD=0.36 $P_k$=0.351 | WD=0.354 $P_k$=0.343 |
| 7 | The Hunger Games | k=46, w=20 | **WD=0.374** $P_k$=0.359 | WD=0.411 $P_k$=0.4 | WD=0.411 $P_k$=0.394 |
| 8 | The Hunger Games: Catching Fire | k=57, w=40 | **WD=0.32** $P_k$=0.305 | WD=0.334 $P_k$=0.324 | WD=0.324 $P_k$=0.327 |
| 9 | The Hunger Games: Mockingjay - Part 1 | k=13, w=97 | WD=0.354 $P_k$=0.35 | **WD=0.335** $P_k$=0.331 | WD=0.349 $P_k$=0.347 |
| 10 | The Hunger Games: Mockingjay - Part 2 | k=21, w=97 | WD=0.326 $P_k$=0.322 | **WD=0.315** $P_k$=0.312 | WD=0.329 $P_k$=0.326 |
| 11 | The Giver | k=,82 w=28 | WD=0.427 $P_k$=0.388 | **WD=0.408** $P_k$=0.372 | WD=0.442 $P_k$=0.414 |
| 12 | The Maze Runner | k=43, w=27 | **WD=0.357 $P_k$=0.338** | WD=0.373 $P_k$=0.36 | WD=0.38 $P_k$=0.367 |
| 13 | Ready Player One | k=95, w=34 | WD=0.365 $P_k$=0.346 | **WD=0.354** $P_k$=0.339 | WD=0.378 $P_k$=0.369 |

# Subtopic boundaries example

- novel "Catching Fire" by S. Collins
- k=11, w=44, $WD_{we}$=1.0, $WD_{we\_oie}$=1.0



Subtopic boundaries (novel)

# Subtopic boundaries example

- subtitles to the film "Ninety Eighty-Four" (1956)
- k=39, w=22, $WD_{we}$=0.423, $WD_{we\_oie}$=0.389



Subtopic boundaries (subtitle)