# Text Technologies for Data Science

## INFR11145

# Indexing
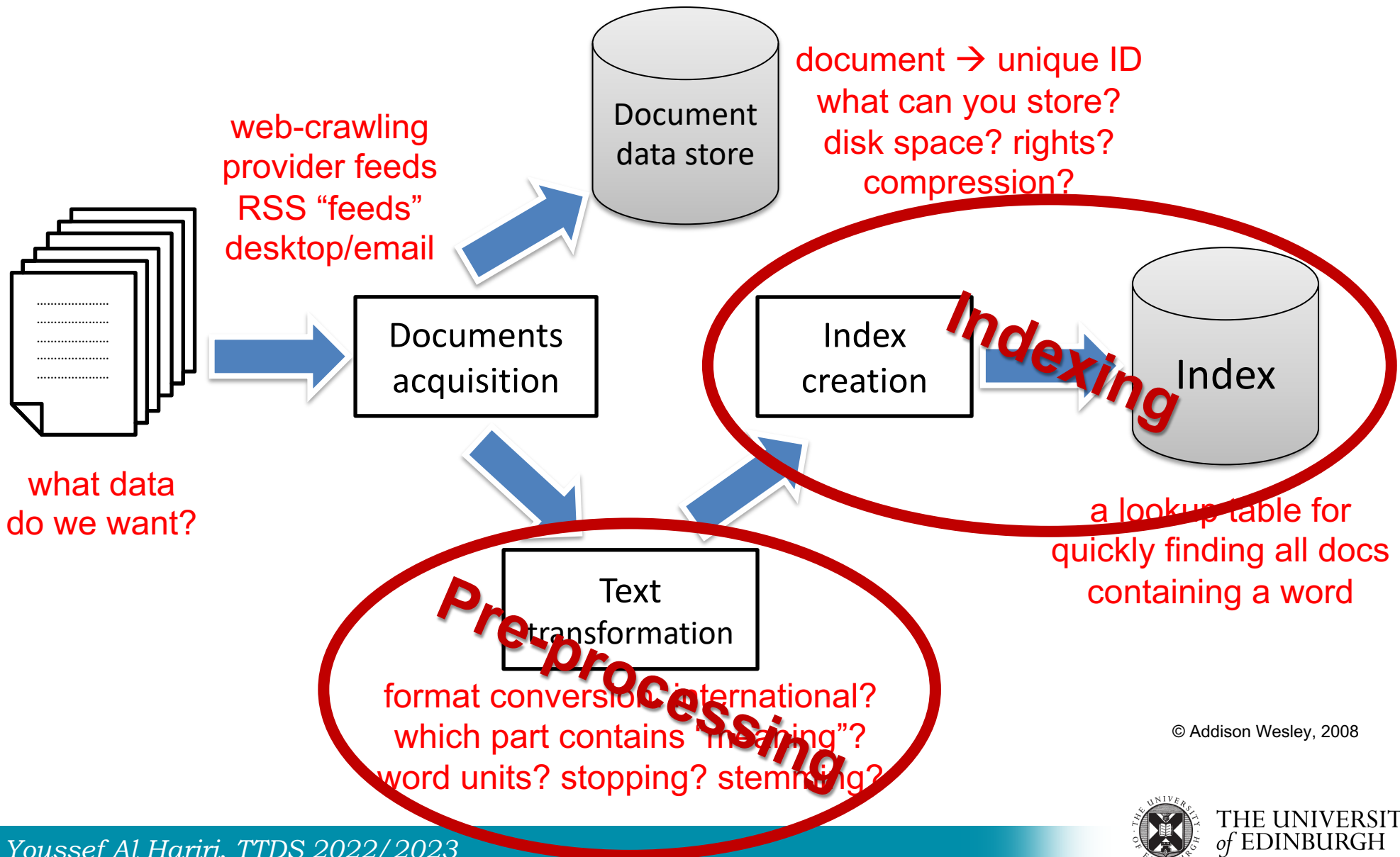
Instructor:
**Youssef Al Hariri**

# Pre-lecture

- Lectures 1-4 → warmup
  - Now, it is getting more serious!

- Lab 1 → slightly different results
  - Tokenisation, stopping, stemming

- Today: two lectures on Indexing
  - expect some knowledge in binary numbers "001011101"

- Announcement of CW1 (Friday, October 6$^{th}$)

- Labs → Do it On Time

- Piazza!!!!

THE UNIVERSITY of EDINBURGH

# Lecture Objectives

- <u>Learn</u> about and <u>implement</u>
- Boolean search
- Inverted index
- Positional index

THE UNIVERSITY *of* EDINBURGH

# Indexing Process



web-crawling
provider feeds
RSS "feeds"
desktop/email

Document data store

document → unique ID
what can you store?
disk space? rights?
compression?

what data
do we want?

Documents acquisition

Index creation

**Indexing**

Index

a lookup table for
quickly finding all docs
containing a word

**Pre-processing**

Text transformation

format conversion? international?
which part contains "meaning"?
word units? stopping? stemming?

© Addison Wesley, 2008

THE UNIVERSITY of EDINBURGH

# Pre-processing output

This is an **exampl**e **sentenc**e of how the **pre-process**ing is **appli**ed to **text** in **inform**ation **retriev**al. It **includ**es: **Token**ization, **Stop Word**s **Remov**al, and **Stem**ming

exampl sentenc pre process appli text inform retriev includ token stop word remov stem

- Add processed terms to the index

- What is "index"?

# Index

- How to match your term in non-linear time?

- Find/Grep:
  Sequential search for term

- Index:
  Find term locations immediately

THE UNIVERSITY *of* EDINBURGH

# Book Index

## Index

THE UNIVERSITY *of* EDINBURGH

# Indexing

- Search engines vs PDF find or grep?

  - Infeasible to scan large collection of text for every "search"

  - Find section that has: "UK and Scotland and Money"?!

- Book Index

  - For each word, list of "relevant" pages

  - Find topic in sub-linear time

- IR Index:

  - Data structure for fast finding terms

  - Additional optimisations could be applied

THE UNIVERSITY *of* EDINBURGH

# Document Vectors

- Represent documents as vectors
  - Vector → document, cell → term
  - Values: term frequency or binary (0/1)
  - All documents → collection matrix

number of occurrence of
a term in a document

# Inverted Index

- Represent terms as vectors
  - Vector → term, cell → document
  - Transpose of the collection matrix
  - Vector: inverted list

| he | drink | ink | likes | pink | thing | wink | |
|---|---|---|---|---|---|---|---|
| 2 | 1 | 0 | 2 | 0 | 0 | 1 | ← **D1:** He likes to wink, he likes to drink |
| 1 | 3 | 0 | 1 | 0 | 0 | 0 | ← **D2:** He likes to drink, and drink, and drink |
| 1 | 1 | 1 | 1 | 0 | 1 | 0 | ← **D3:** The thing he likes to drink is ink |
| 1 | 1 | 1 | 1 | 1 | 0 | 0 | ← **D4:** The ink he likes to drink is pink |
| 1 | 1 | 1 | 1 | 1 | 0 | 1 | ← **D5:** He likes to wink, and drink pink ink |

THE UNIVERSITY *of* EDINBURGH

# Boolean Search

- Boolean: exist / not-exist

- Multiword search: logical operators (AND, OR, NOT)

- Example
  - Collection: search Shakespeare's Collected Works
  - Boolean query: Brutus AND Caesar AND NOT Calpurnia

- Build a **Term-Document Incidence Matrix**
  - Which term appears in which document
  - Rows are terms
  - Columns are documents

# Collection Matrix

**Documents**

| | Antony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth |
|---|---|---|---|---|---|---|
| **Antony** | 1 | 1 | 0 | 0 | 0 | 1 |
| **Brutus** | 1 | 1 | 0 | 1 | 0 | 0 |
| **Caesar** | 1 | 1 | 0 | 1 | 1 | 1 |
| **Calpurnia** | 0 | 1 | 0 | 0 | 0 | 0 |
| **Cleopatra** | 1 | 0 | 0 | 0 | 0 | 0 |
| **mercy** | 1 | 0 | 1 | 1 | 1 | 1 |
| **worser** | 1 | 0 | 1 | 1 | 1 | 0 |

**Terms**

1 if *document* contains *term*, 0 otherwise

Query: Brutus AND Caesar AND NOT Calpurnia
Apply on rows: **110100** AND **110111** AND !(**010000**) = **100100**

# Bigger collections?

- Consider *N* = 1 million documents, each with about 1000 words.

- *n* = 1M x 1K = 1B words
  ➜ Heap's law ➜ *v* ≈ 500K

- Matrix size = 500K unique terms x 1M documents = 0.5 trillion 0's and 1's entries!

- If all words appear in many documents
  ➜ max{count(1's)} = *N* * doc. length = 1B

- Actually, from Zip's law ➜ 250k terms appears once!

- Collection matrix is extremely **sparse**. (*mostly 0's*)

term$_x$

1
0
0
0
1
1
0
0
0
0
0
0
0
0
0
0
0
0
0
0
0
0
1
0
0
0
0
0
0
0
0
1
.
.
.
.
.

1M

THE UNIVERSITY *of* EDINBURGH

# Inverted Index: Sparse representation

- For each term $t$, we must store a list of all documents that contain $t$.
  - Identify each by a **docID**, a document serial number

| **Posting** |

| **Brutus** | → | 1 | 2 | 4 | 11 | 31 | 45 | 173 |

| **Caesar** | → | 1 | 2 | 4 | 5 | 6 | 16 | 57 | 132 |

| **Calpurnia** | → | 2 | 31 | 54 | 101 |

**Dictionary**

**Postings List**

**Doc number (sorted)**

THE UNIVERSITY *of* EDINBURGH

# Inverted Index Construction

**Documents to be indexed**

Friends, Romans, countrymen

●
●
●

*Tokenizer*

**Token stream**

| Friends | Romans | Countrymen |

*Normaliser*

**Terms (modified tokens)**

| friend | roman | countryman |

*Indexer*

**Inverted index**

| *friend* | → | 2 | → | 4 | → |
| *roman* | → | 1 | → | 2 | → |
| *countryman* | → | 3 | → | 9 | → |

# Step 1: Term Sequence

**Doc 1**

I did enact Julius Caesar I was
killed i' the Capitol;  Brutus killed
me.

**Doc 2**

So let it be with Caesar. The
noble Brutus hath told you
Caesar was ambitious

**Sequence of
(term, Doc ID) pairs** →

| Term | docID |
|---|---|
| I | 1 |
| did | 1 |
| enact | 1 |
| julius | 1 |
| caesar | 1 |
| I | 1 |
| was | 1 |
| killed | 1 |
| i' | 1 |
| the | 1 |
| capitol | 1 |
| brutus | 1 |
| killed | 1 |
| me | 1 |
| so | 2 |
| let | 2 |
| it | 2 |
| be | 2 |
| with | 2 |
| caesar | 2 |
| the | 2 |
| noble | 2 |
| brutus | 2 |
| hath | 2 |
| told | 2 |
| you | 2 |
| caesar | 2 |
| was | 2 |
| ambitious | 2 |

# Step 2: Sorting

- ## Sort by:

### 1) Term

**then**

### 2) Doc ID

| Term | docID |
|------|-------|
| I | 1 |
| did | 1 |
| enact | 1 |
| julius | 1 |
| caesar | 1 |
| I | 1 |
| was | 1 |
| killed | 1 |
| i' | 1 |
| the | 1 |
| capitol | 1 |
| brutus | 1 |
| killed | 1 |
| me | 1 |
| so | 2 |
| let | 2 |
| it | 2 |
| be | 2 |
| with | 2 |
| caesar | 2 |
| the | 2 |
| noble | 2 |
| brutus | 2 |
| hath | 2 |
| told | 2 |
| you | 2 |
| caesar | 2 |
| was | 2 |
| ambitious | 2 |

**Sorting** →

| Term | docID |
|------|-------|
| ambitious | 2 |
| be | 2 |
| brutus | 1 |
| brutus | 2 |
| capitol | 1 |
| caesar | 1 |
| caesar | 2 |
| caesar | 2 |
| did | 1 |
| enact | 1 |
| hath | 1 |
| I | 1 |
| I | 1 |
| i' | 1 |
| it | 2 |
| julius | 1 |
| killed | 1 |
| killed | 1 |
| let | 2 |
| me | 1 |
| noble | 2 |
| so | 2 |
| the | 1 |
| the | 2 |
| told | 2 |
| you | 2 |
| was | 1 |
| was | 2 |
| with | 2 |

# Step 3: Posting

1. Multiple term entries in a single document are merged

2. Split into Dictionary and Postings

3. Doc. Frequency (*df*) information is added

| Term | docID |
|---|---|
| ambitious | 2 |
| be | 2 |
| brutus | 1 |
| brutus | 2 |
| capitol | 1 |
| caesar | 1 |
| caesar | 2 |
| caesar | 2 |
| did | 1 |
| enact | 1 |
| hath | 1 |
| I | 1 |
| I | 1 |
| i' | 1 |
| it | 2 |
| julius | 1 |
| killed | 1 |
| killed | 1 |
| let | 2 |
| me | 1 |
| noble | 2 |
| so | 2 |
| the | 1 |
| the | 2 |
| told | 2 |
| you | 2 |
| was | 1 |
| was | 2 |
| with | 2 |

| term | doc. freq. | → | postings lists |
|---|---|---|---|
| ambitious | 1 | → | 2 |
| be | 1 | → | 2 |
| brutus | 2 | → | 1 → 2 |
| capitol | 1 | → | 1 |
| caesar | 2 | → | 1 → 2 |
| did | 1 | → | 1 |
| enact | 1 | → | 1 |
| hath | 1 | → | 2 |
| i | 1 | → | 1 |
| i' | 1 | → | 1 |
| it | 1 | → | 2 |
| julius | 1 | → | 1 |
| killed | 1 | → | 1 |
| let | 1 | → | 2 |
| me | 1 | → | 1 |
| noble | 1 | → | 2 |
| so | 1 | → | 2 |
| the | 2 | → | 1 → 2 |
| told | 1 | → | 2 |
| you | 1 | → | 2 |
| was | 2 | → | 1 → 2 |
| with | 1 | → | 2 |

THE UNIVERSITY *of* EDINBURGH

# Inverted Index: matrix → postings

| he | drink | ink | likes | pink | thing | wink | |
|----|-------|-----|-------|------|-------|------|---|
| 2 | 1 | 0 | 2 | 0 | 0 | 1 | ← **D1:** He likes to wink, he likes to drink |
| 1 | 3 | 0 | 1 | 0 | 0 | 0 | ← **D2:** He likes to drink, and drink, and drink |
| 1 | 1 | 1 | 1 | 0 | 1 | 0 | ← **D3:** The thing he likes to drink is ink |
| 1 | 1 | 1 | 1 | 1 | 0 | 0 | ← **D4:** The ink he likes to drink is pink |
| 1 | 1 | 1 | 1 | 1 | 0 | 1 | ← **D5:** He likes to wink, and drink pink ink |

*he* → 1 2 3 4 5

*drink* → 1 2 3 4 5

*ink* → 3 4 5

*pink* → 4 5

*thing* → 3

*wink* → 1 5

THE UNIVERSITY
*of* EDINBURGH

# Inverted Index: with frequency

- Boolean: term → DocIDs list

- Frequency: term → touples (DocID,count(term)) lists

| | | | | | |
|---|---|---|---|---|---|
| *he* → | 1:2 | 2:1 | 3:1 | 4:1 | 5:1 |
| *drink* → | 1:1 | 2:3 | 3:1 | 4:1 | 5:1 |
| *ink* → | 3:1 | 4:1 | 5:1 | | |
| *pink* → | 4:1 | 5:1 | | | |
| *thing* → | 3:1 | | | | |
| *wink* → | 1:1 | 5:1 | | | |

appeared in D2 3 times

THE UNIVERSITY *of* EDINBURGH

# Query Processing

- Find documents matching query {ink **AND** wink}
    1. Load inverted lists for each query word
    2. Merge two postings lists → **Linear merge**

- Linear merge → O($n$)
  $n$: total number of posts for all query words



ink → | 3:1 | 4:1 | 5:1 |

wink → | 1:1 | 5:1 |

Matches

1: f(0,1)

3: f(1,0)

4: f(1,0)

5: f(1,1)

THE UNIVERSITY
*of* EDINBURGH

# Phrase Search

- Find documents matching query "pink ink"
    1. Find document containing both words
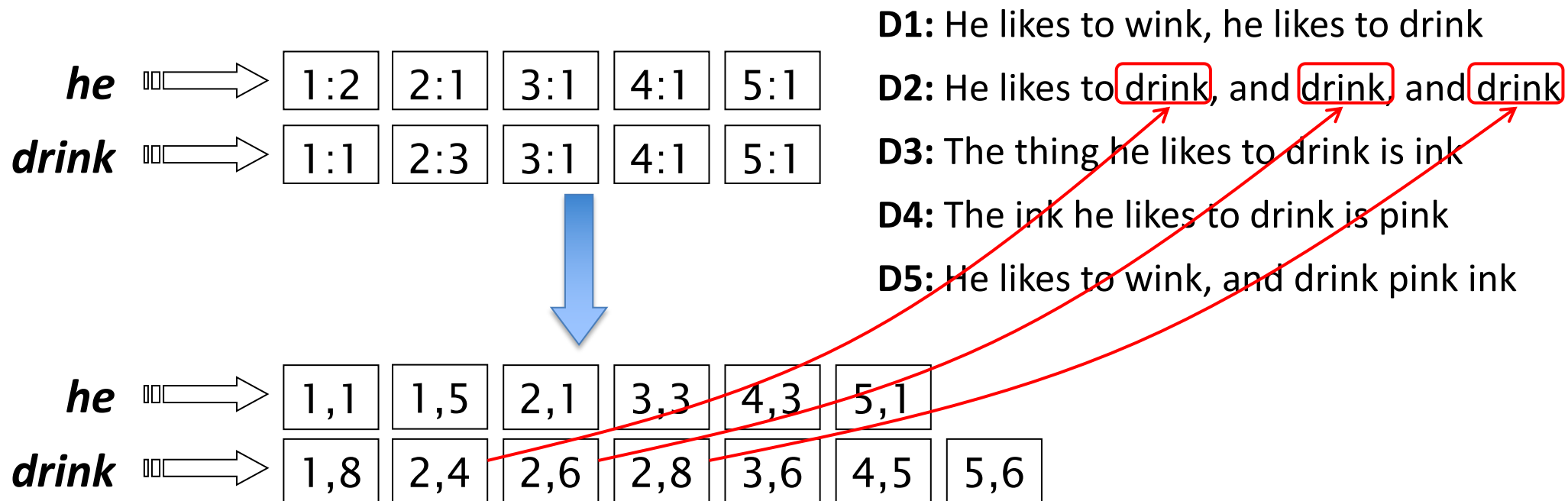    2. Both words has to be a phrase

- Bi-gram Index:

He likes to wink, and drink pink ink  →  Convert to bigrams

He_likes likes_to to_wink wink_and and_drink drink_pink pink_ink

- Bi-gram Index, issues:
    - Fast, but index size will explode!
    - What about trigram phrases?
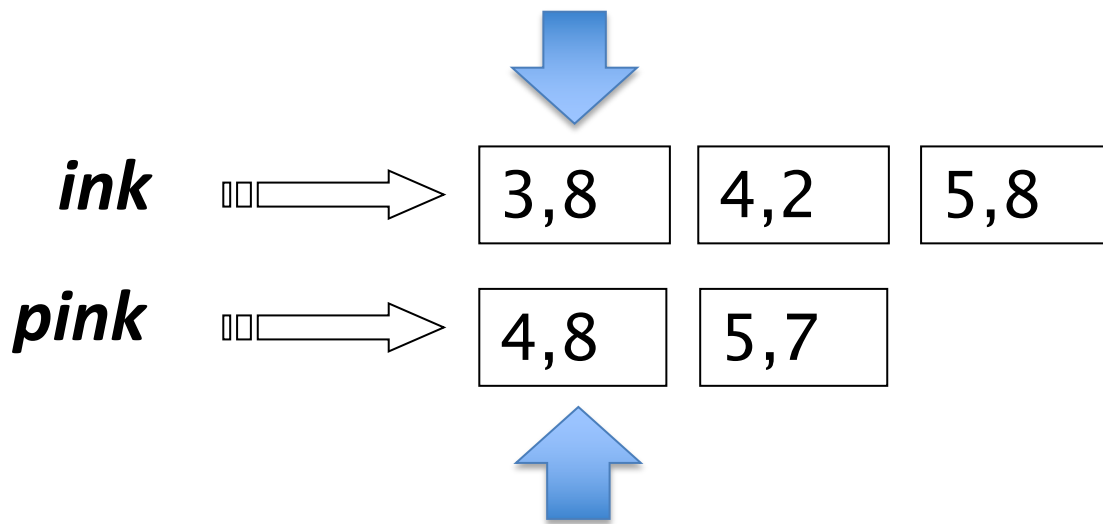    - What about proximity? "ink is pink"

# Proximity Index

- Terms positions is embedded to the inv. Index
    - Called proximity/positional index
    - Enables phrase and proximity search
    - Toubles (DocID, term position)

**he** ⇒ | 1:2 | 2:1 | 3:1 | 4:1 | 5:1 |

**drink** ⇒ | 1:1 | 2:3 | 3:1 | 4:1 | 5:1 |

**D1:** He likes to wink, he likes to drink

**D2:** He likes to drink, and drink, and drink

**D3:** The thing he likes to drink is ink

**D4:** The ink he likes to drink is pink

**D5:** He likes to wink, and drink pink ink

**he** ⇒ | 1,1 | 1,5 | 2,1 | 3,3 | 4,3 | 5,1 |

**drink** ⇒ | 1,8 | 2,4 | 2,6 | 2,8 | 3,6 | 4,5 | 5,6 |

THE UNIVERSITY of EDINBURGH

# Query Processing: Proximity

- Find documents matching query "pink ink"
  1. Use **Linear merge**
  2. Additional step: check terms positions

- **Proximity search**:
  
  *pos(term1) – pos(term2) < |w|* → #5(pink,ink)

**Matches**

ink  →  | 3,8 | | 4,2 | | 5,8 |

pink  →  | 4,8 | | 5,7 |

3: f(1,0) = **0**

4: f(1,1) = ? =
pos(ink) – pos(pink) == 1?

5: f(1,1) = ? =
pos(ink) – pos(pink) == 1?

THE UNIVERSITY
*of* EDINBURGH

# Proximity search: data structure

- Possible data structure:

  <term: df;

   DocNo: pos1, pos2, pos3

   DocNo: pos1, pos2, pos3

   ……. >

- Example:

  <*be*: 993427;

   *1*: 7, 18, 33, 72, 86, 231;

   *2*: 3, 149;

   *4*: 17, 191, 291, 430, 434;

   *5*: 363, 367, …>

THE UNIVERSITY *of* EDINBURGH

# Practical

# Summary

- Document Vector

- Term Vector

- Inverted Index

- Collection Matrix

- Posting

- Proximity Index

- Query Processing → Linear merge

THE UNIVERSITY *of* EDINBURGH

# Resources

- Textbook 1: Intro to IR, Chapter 1 & 2.4

- Textbook 2: IR in Practice, Chapter 5

- Lab 2

THE UNIVERSITY
*of* EDINBURGH