# A New Corpus for the Evaluation of Arabic Intrinsic Plagiarism Detection

Imene Bensalem[1], Paolo Rosso[2], Salim Chikhi[1]

[1]Constantine 2 University, Algeria
[2]Universitat Politècnica de Valeència, Spain
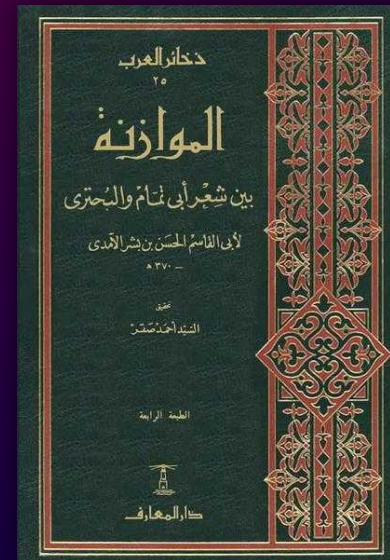
# Outline


كلمات أفكار
Ideas
words

# Outline

- Arabic Text Reuse

- Plagiarism in Arab World

- Plagiarism Detection Approaches
  - Intrinsic Plagiarism Detection

- Arabic Plagiarism Detection

- Building Plagiarism Detection Corpora

- InAra: the First Corpus for the Evaluation of Arabic Intrinsic Plagiarism Detection

# Outline

# Arabic Text Reuse

☐ Text reuse studies were an active field in Arabic literature criticism in the middle ages (Al Manasrah 2009)

☐ Text reuse was spread and studied especially in the context of poetry

Medieval books on plagiarism in poetry

# Arabic Text Reuse

□ All kinds of text reuse in literature were called "literary theft" by the Arab medieval critics,

□ BUT, "literary theft" was classified into many types and levels, some of which were acceptable or even considered an art that could be done only by the very prominent poets! (Tabana 1956)



Words used to express the different kinds of literary thefts

# Arabic Text Reuse

Till now, the word "Plagiarism" does not have a unique term in Arabic.
Some of the used terms: literary theft , scientific theft, arrogation…

But there is a tendency of using the word انتحال /intihal/ that means arrogation of authorship (in poetry)
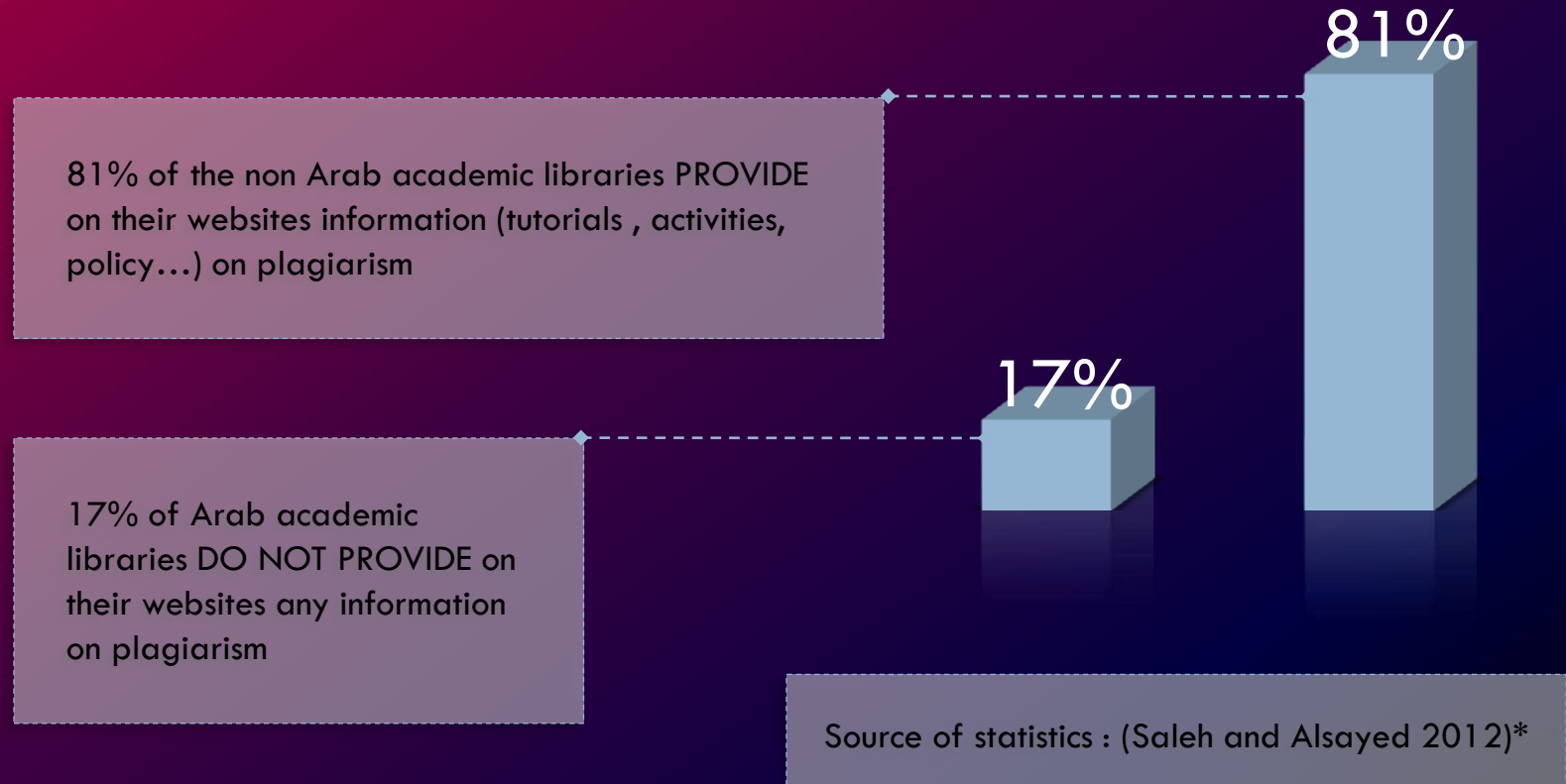
انتحال

# Plagiarism in Arab World

## Study (Al-Jundy 2013)*

- 18 magister candidates from different universities in Egypt were asked to write 3 essays : 2 in Arabic and 1 in English on topics related to library sciences. It was reported that:

- All of them plagiarized

- The percentage of plagiarism in each plagiarized essay was between 45%-76% , it reached 81% in essays written in English

- They forgot or do not know the importance of giving acknowledge in academic writing

- They plagiarized despite the warning on plagiarism because they did not know about the existence of any plagiarism detection software!

# Plagiarism in Arab World

One of plagiarism reasons in Arab world

81% of the non Arab academic libraries PROVIDE on their websites information (tutorials , activities, policy...) on plagiarism

**81%**

**17%**

17% of Arab academic libraries DO NOT PROVIDE on their websites any information on plagiarism

Source of statistics : (Saleh and Alsayed 2012)*

*عماد عيسى صالح، أماني محمد السيد. دور المكتبات الأكاديمية في منع السرقات العلمية واكتشافها: دراسة استكشافية لخدمات المكتبات وبرمجيات كشف الانتحال. المؤتمر الدولي للتعلم الالكتروني في الوطن العربي، القاهرة 9-11 يوليو 2012.

# To Follow Plagiarism Cases…

السرقات الأدبية@ Facebook

List of plagiarism incidents @ Wikipedia

Tweeter@سرقوني#

http://bader59.com

# Outline



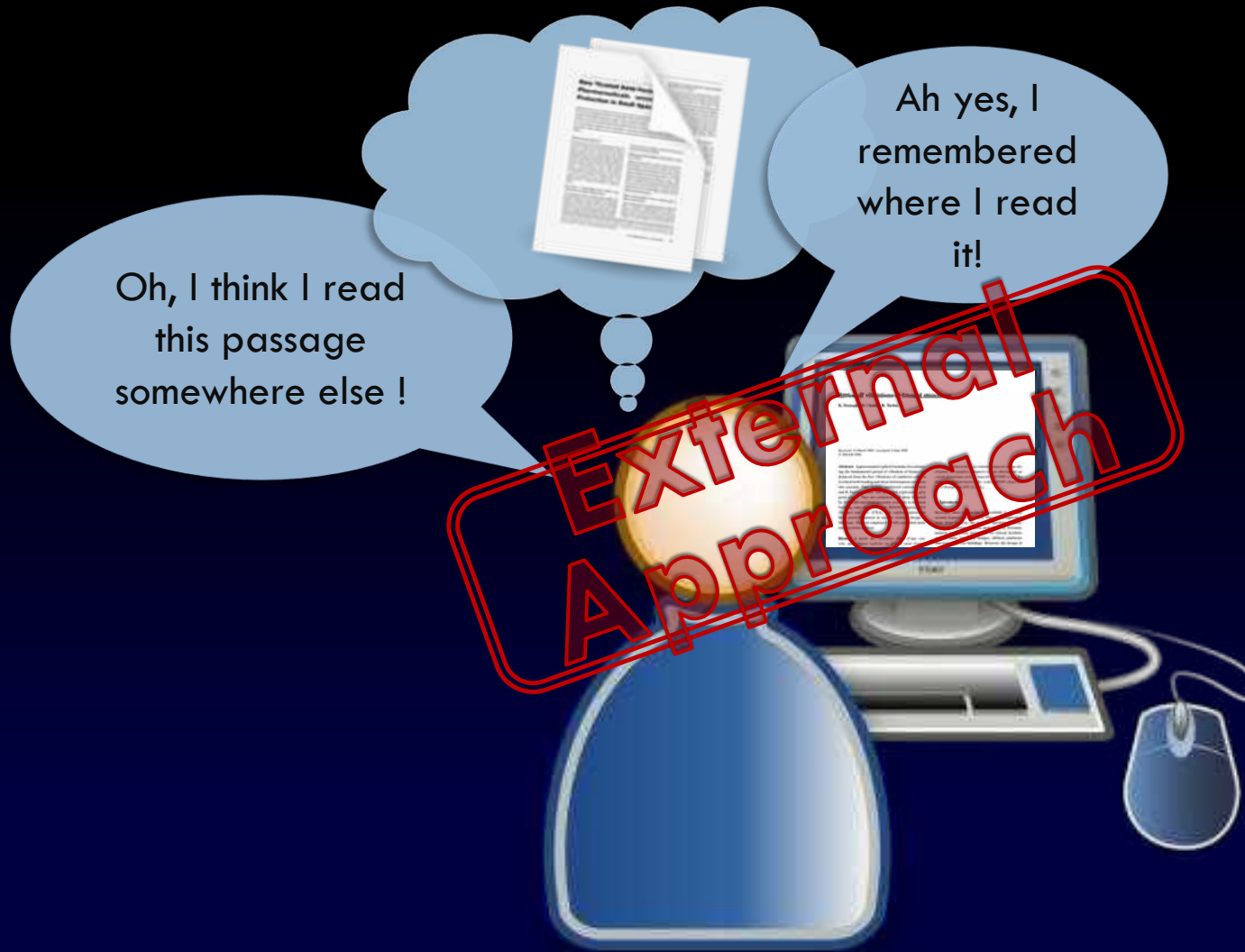- Arabic Text Reuse
- Plagiarism in Arab World

# Outline



□ Arabic Text Reuse

□ Plagiarism in Arab World

□ **Plagiarism Detection Approaches**

   □ Intrinsic Plagiarism Detection

□ **Arabic Plagiarism Detection**

□ Building Plagiarism Detection Corpora

□ InAra: the First Corpus for the Evaluation Arabic Intrinsic Plagiarism Detection

# Plagiarism Detection Approaches

# Plagiarism Detection Approaches

Intrinsic Plagiarism Detection Techniques

Intrinsic plagiarism detection methods building blocks (after (Stein et al. 2011)) with examples from the state-of-the-art techniques

**Pre-processing**
- Plagiarism-free document (Stamatatos 2011)

**Segmentation**
- Topical (Muhr et al. 2010)
- Sentence-based (Zechner et al. 2009)
- Character-based (Stamatataros 2009)
- Sliding window
- Uniform length (Stein et al. 2011)
- Natural chunks (Meyer zu Eissen et al. 2006, 2007)

**Quantifying the Writing style**
- Complexity (Seaward and Matwin 2009)
- Syntactic features
- Vocabulary richness (Meyer zu Eissen et al. 2006,)
- Lexical features
- N-gram (Stamatatos 2009)
- Character features

**Outliers detection**
- Density-based method (Stein et al. 2011)
- Similarity/difference-based method (Stamatatos 2009)

**Post-processing**
- Concatenation of the adjacent outliers (Kestemont et al. 2011)

# Intrinsic Plagiarism Detection Corpora

Meyer zu Eissen et al. 2006 Corpus

More than 450 suspicious documents
Computer science papers from ACM

Stein and Meyer zu Eissen 2007 Corpus

Scientific German documents from German universities
1600 non-plagiarized sections + 1500 plagiarized sections.

**PAN 2009 Corpus**

Separate corpora for the external and the intrinsic tasks
The sub-corpus for the intrinsic task: 3091 suspicious documents
Free books from Gutenberg

**PAN 2010 Corpus**

One corpus for the external and the intrinsic task :
30% of documents are used for the intrinsic task
Free books from Gutenberg

**PAN 2011 Corpus**

Separate corpora for the external and the intrinsic tasks
The sub-corpus for the intrinsic task: 4753 suspicious documents
Free books from Gutenberg

**PAN 2012 Corpus**

Authorship clustering corpus
Less than 10 documents composed of segments from many authors

# Arabic Plagiarism Detection

- Few works evaluated on different corpora
- All of which are in external plagiarism detection

Cloud of some keywords from Arabic plagiarism detection papers

# Arabic Plagiarism Detection Software



Source: Alzahrany papers

# Arabic Plagiarism Detection Software



http://qarnet.com

# Intrinsic Plagiarism Detection
## State and Challenges

**In general**

- It is still in its infancy
- It is a reduction of the authorship verification task (one-class classification problem) BUT
  - The class is not well defined (text from which the author writing style is extracted is noisy with the plagiarized chunks)
  - Verify the authorship of short segment of text

**In Arabic**

- Very few works on Arabic Stylometry and Authorship Analysis (the source of techniques)
- May need religion quotation detection : Arabic text may contain quotation from religious text

# Outline

- Arabic Text Reuse
- Plagiarism in Arab World

- Plagiarism Detection Approaches
  - Intrinsic Plagiarism Detection
- Arabic Plagiarism Detection

- Building Plagiarism Detection Corpora
- InAra: the First Corpus for the Evaluation of Arabic Intrinsic Plagiarism Detection

كلمات أفكار
Ideas
words

# Outline

- Arabic Text Reuse
- Plagiarism in Arab World
- Plagiarism Detection Approaches
  - Intrinsic Plagiarism Detection
  - Arabic Plagiarism Detection

- Building Plagiarism Detection Corpora
- InAra: the First Corpus for the Evaluation of Arabic Intrinsic Plagiarism Detection

# Building Plagiarism Detection Corpora
## Automatic Approach



Plagiarism metadata

Suspicious document

topic
=

author
≠

**2** Select text passages randomly form source documents and Insert them in random positions in the target document

**1** Text Document Compilation

Target documents

Source documents

# Building Plagiarism Detection Corpora
## Manual Approach
## (Crowdsoursing)

# Building Plagiarism Detection Corpora
## Manual Approach

Source Documents

Write an essay on Topic X

PONY

# Building Plagiarism Detection Corpora
## Manual Approach

# Building Plagiarism Detection Corpora
## Manual Approach

# Building Plagiarism Detection Corpora
## Manual Approach

# Building Plagiarism Detection Corpora
## Manual Approach

# Building Plagiarism Detection Corpora
## Manual Approach

# Building Plagiarism Detection Corpora
## Comparison of Approaches

|  | Automatic approach | Manual approach |
|---|:---:|:---:|
| Document length variety | ✔ | ✘ |
| Real plagiarism scenario ? | ✘ | ✔ |
| Copyright issue on suspicious documents | ✘ | ✔ |
| Cost (Material and human resource) | ✔ | ✘ |

# InAra: the First Corpus for the Evaluation Arabic Intrinsic Plagiarism Detection

- ☐ Built using the automatic approach.

- ☐ Target and source document were collected from Wikisource.

- ☐ Copyright-free documents.

- ☐ Documents tagged with topics and author name.

**WIKISOURCE**

| | |
|---|---|
| **English** The Free Library 358,000+ pages | **Français** La bibliothèque libre 84,000+ pages |
| **Русский** Свободная библиотека 196,000+ статей | **Español** La biblioteca libre 51,000+ páginas |
| **Deutsch** Die freie Quellensammlung 83,000+ Seiten | **Italiano** La biblioteca libera 62,000+ pagine |
| **Português** A biblioteca livre 87,000+ páginas | **עברית** הספריה החופשית 33,000+ מאמרים |
| **Polski** wolna biblioteka 32,000+ strony | **العربية** المكتبة الحرة 20,000+ صفحة |

# Criteria of Target Documents

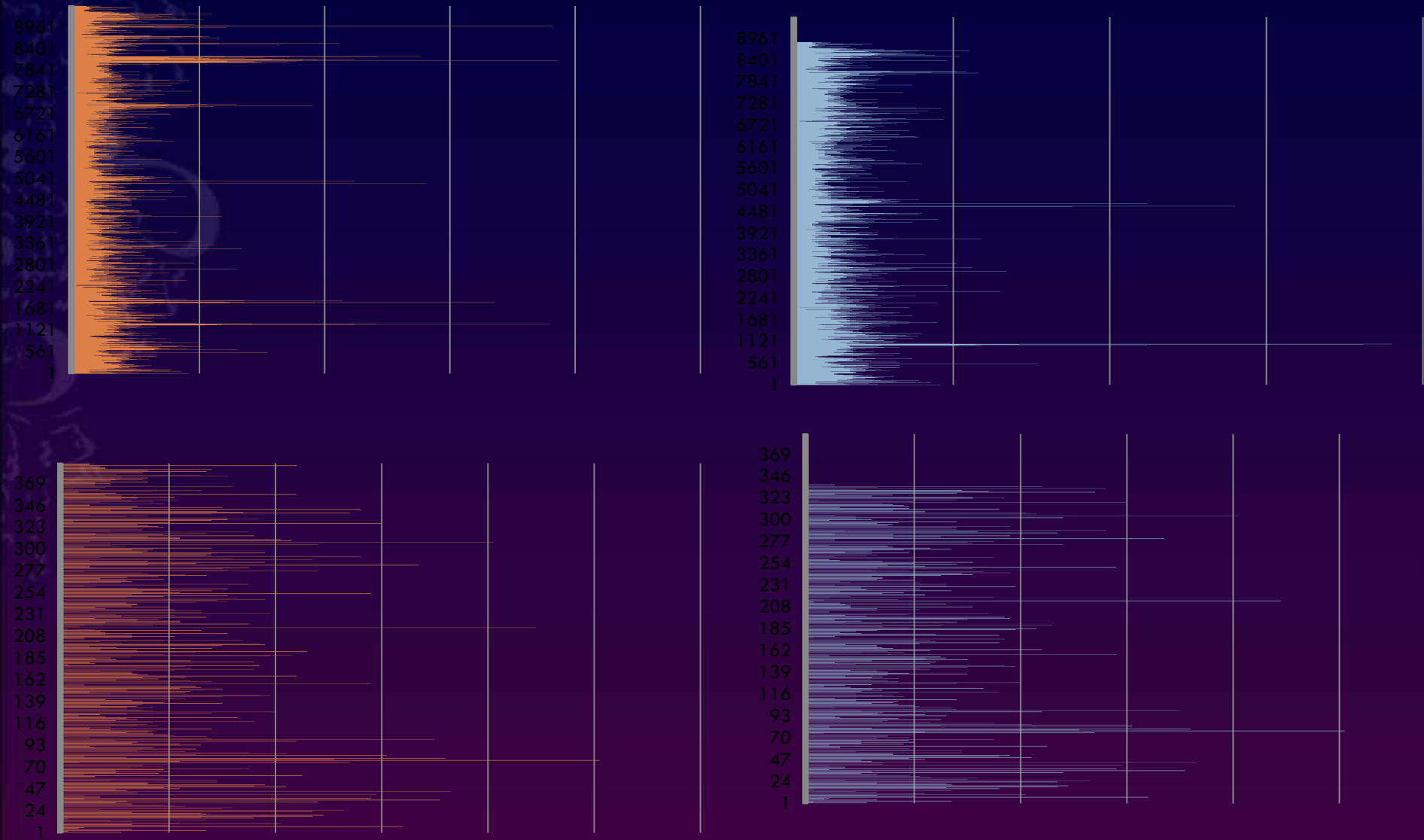| | |
|---|---|
| Criteria to not complicate the detection task further more | □ Written by one author only: ~~Wikipedia~~<br><br>□ Should not include much of text reuse or many quotations: ~~News stories~~, ~~religion books with many quotations~~ |
| Criteria to obtain reliable writing style model | □ Should not be too short: ~~News stories~~<br><br>□ Should be well edited: ~~texts without punctuation~~ |

# Problem Encountered when Building the Corpus

- Limited copyright-free source of text

- Poor quality of the online free texts: not well punctuated and edited

- A high percentage from the free text are books on religious topics => contain many quotations => not very suited for a corpus to evaluate the intrinsic approach

# Why a Corpus in Arabic ?

- The ability to test language-dependent methods i.e. that take into account Arabic peculiarities (e.g. diacritics) or based on processing like parsing, stemming…etc

- Writing style discriminators in English are not necessarily discriminator in Arabic, e.g. average word length (Bensalem et al. 2012) (The next slide for other experiment on the sentence length)

- The importance of corpora to foster research

# Visualization of Sentence length of some documents before and after the insertion of plagiarism (Sentence length does not seem to be a good discriminator in Arabic text

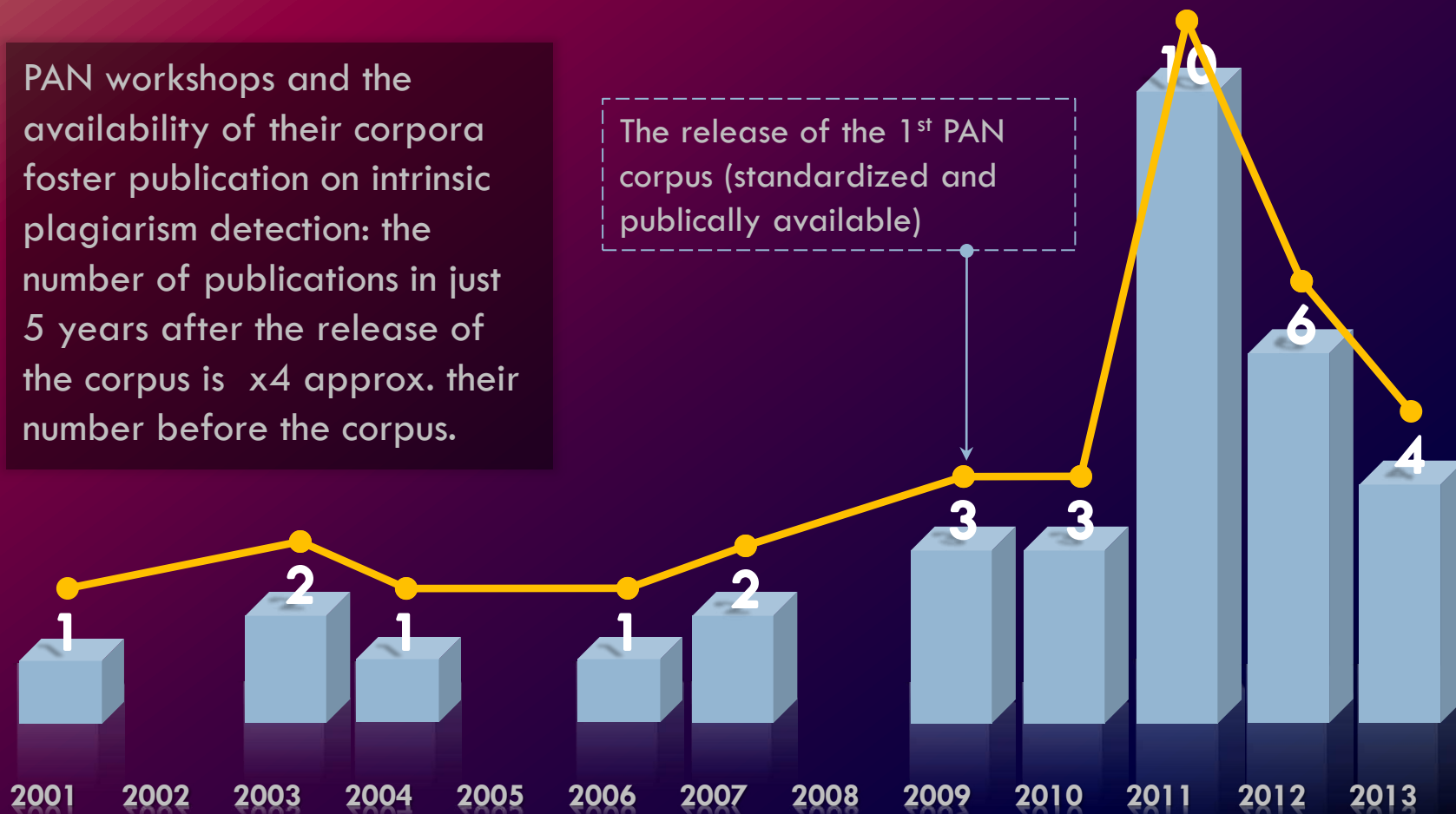# Why a Corpus in Arabic ?

- ☐ The ability to test language-dependent methods i.e. that take into account Arabic peculiarities (e.g. diacritics) or based on processing like parsing, stemming…etc

- ☐ Writing style discriminators in English are not necessarily discriminator in Arabic, e.g. average word length (Bensalem et al. 2012)

- ☐ The importance of the corpus to foster research (The next slide)

# Number of Publications on Intrinsic Plagiarism Detection from 2001 to 2013

PAN workshops and the availability of their corpora foster publication on intrinsic plagiarism detection: the number of publications in just 5 years after the release of the corpus is x4 approx. their number before the corpus.

The release of the 1st PAN corpus (standardized and publically available)

1   2   1   1   2   3   3   10   6   4

2001  2002  2003  2004  2005  2006  2007  2008  2009  2010  2011  2012  2013

# InAra Corpus Statistics

| Document statistics | | | |
|---|---|---|---|
| **Total number of documents** | 1024 | | |
| **Plagiarism percentage per document** | | **Document length** | |
| Null (0%) | 20% | Very Short (1-3 pages) | 46% |
| Hardly ]0%  10%] | 24% | Short (3-15 pages) | 37% |
| Few ]10%  30%] | 32% | Medium (15-100 pages) | 12% |
| Medium ]30%  60%] | 24% | Long (>100 pages) | 05% |
| **Plagiarism cases statistics** | | | |
| **Total number of plagiarism cases** | 2833 | | |
| **Plagiarism cases length** | | **Number of plagiarism cases per document** | |
| Very short (some sentences) | 09% | Null (0) | 20% |
| Short (some paragraphs) | 40% | Few ]0  5] | 69% |
| Medium (around 1 page) | 21% | Medium ]5  15] | 08% |
| Long (many pages) | 30% | Much ]15  45] | 03% |

# How to Get InAra Corpus ?

# Conclusion

Intrinsic plagiarism detection is a challenging task that need more attention from the research community

We hope InAra corpus will help to foster research in this area in Arabic language