# Patent Document Retrieval and Classification at KAIST

**KAIST CS Dept. / BOLA**

**2005. 12. 7.**

Jae-Ho Kim, Jin-Xia Huang, Ha-Yong Jung, Key-Sun Choi

# *Contents*

❖ Introduction

❖ Our approach for Retrieval and Categorization
  - Component-by-component retrieval
  - kNN-based approach

❖ Retrieval and Categorization System

❖ Experimental Results

❖ Conclusions

# *Introduction*

- ❖ Patent Retrieval Task
  - ▪ Document retrieval subtask
  - ▪ **Theme categorization subtask**
    - • Main topic in today's talk

- ❖ Our assumption
  - ▪ Only some parts of patent document are useful for retrieval and categorization
  - ➔ To find important parts and to utilize them well

# Characteristic of Patent Documents

## ❖ *Patent document is structured*

| **Normative section** | |
|---|---|
| | <DOCNO>PATENT-JA-UPA-1995-000001< |
| **<Bibliography>**<br>[publication date]<br>[title of invention] | <SDO BIJ><br>(43)【公開日】平成7年（1995）1月6日<br>(54)【発明の名称】スラリ散布を行う土壌作業機<br>...... |
| **<Abstract>**<br>[purpose]<br>[composition] | <SDO ABJ><br>【目的】スラリの処理と土壌作業を同時に行うことで、……<br>【構成】トラクタとスラリを積載したバキュムカーとの間に …… |
| **<Claims>**<br>[claim1]<br>[claim2] | <SDO CLJ><br>【請求項1】バキウムカーを牽引<br>【請求項2】トラクタに対して3点リ |
| **<Description>**<br>[industrial application field]<br>[problem to be solved]<br>[means of solving problems]<br>[operation]<br>[embodiment examples]<br>[effects of invention] | <SDO DEJ><br>【産業上の利用分野】本発明はスラリ散布を行う土壌作業機に関し、……<br>【発明が解決しようとする課題】このようなスラリを圃場に供給する……<br>【課題を解決するための手段】上述のような目的を達成するために、……<br>【作用】本発明のスラリ散布を行う土壌作業機は、……<br>【実施例】以下、本発明を採用した土壌作業機について添付した図面に……<br>【発明の効果】以上の説明から明らかなように、…… |
| **<Explanation of Drawings>** | <SDO EDJ>【図1】本発明のスラリ散布を行う土壌作業機の側面図である。 |
| **<Drawings>** | <SDO DRJ>【図1】 |

**Detailed component**

*Applicant-defined tags*

4

# *Usefulness of Detailed Components*

❖ [prior art] and [application field]

- Including much information related to technical background and technical field
  - They can be more helpful to categorize patent documents

❖ [purpose] and [means of solving problems]

- Representing the whole patent document
- Often used in the <Abstract> section

❖ <description> section

- including many noises
  - ➔ Selecting useful detailed components in the section

Detailed components ➔ Major features

# *Our Approach for Categorization*

❖ kNN-based patent categorization

- Retrieving k similar documents from training set
- Classifying a given patent into the theme codes of *k* similar documents

❖ Motivation

- Word-bag Vectors
  - are often used in many machine learning method
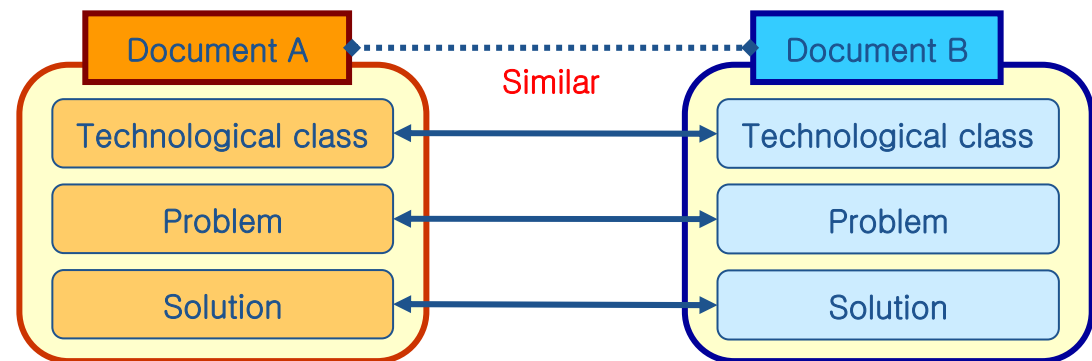  - too large vocabulary size in a large-scaled training set

<div align="center">

**kNN-based approach**

</div>

❖ Component-by-component comparison

- If two documents are in the same technical classes and have the same problem and solution (method)

  ➔ They are similar

| Document A | | Document B |
|---|---|---|
| Technological class | ⬌ Similar | Technological class |
| Problem | ⬌ | Problem |
| Solution | ⬌ | Solution |

❖ Motivation

- Document–by-document comparison
  - Mixed information may cause error

## Component-by-component comparison

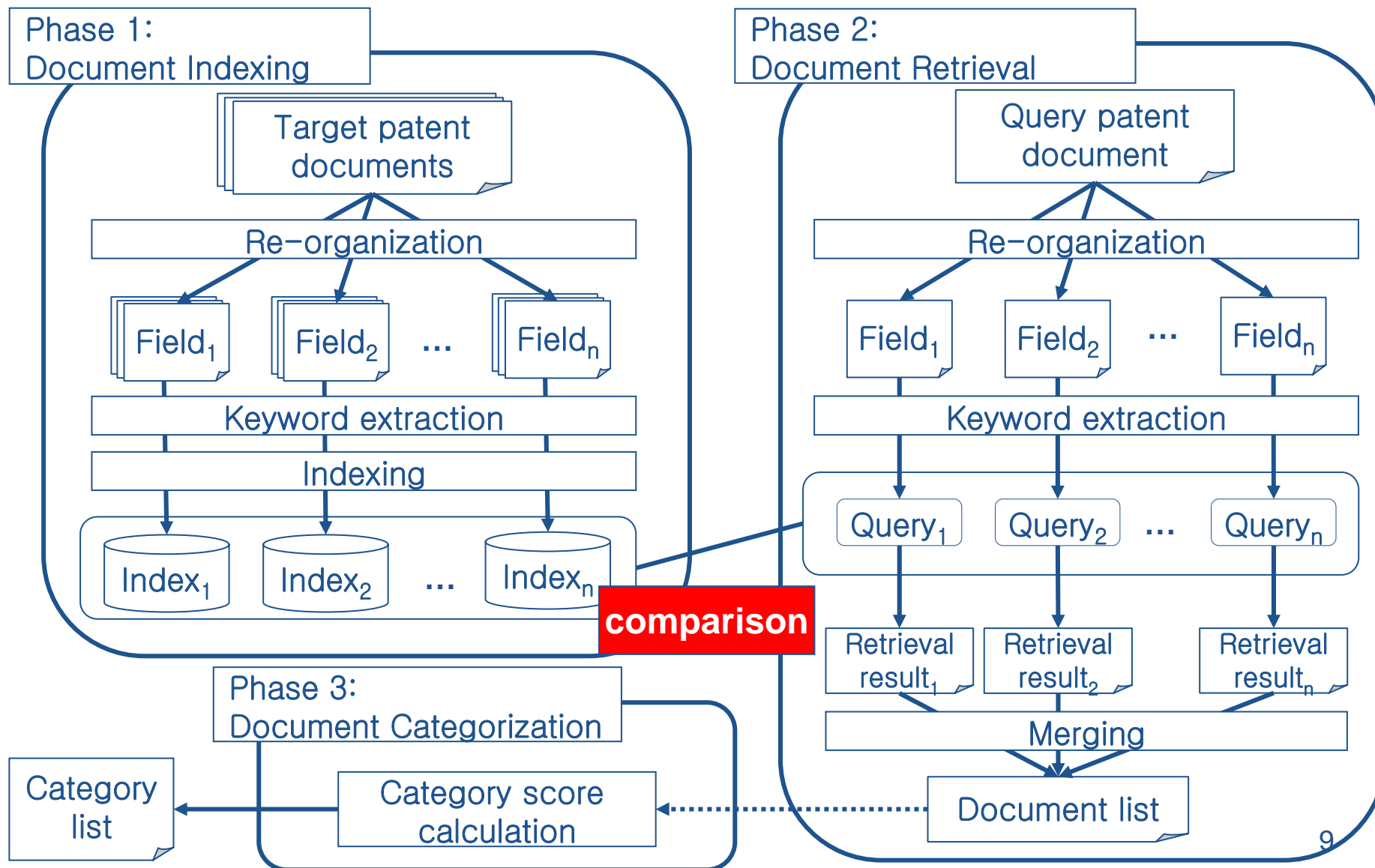# Re-organization for component-by-component comparison

<DOCNO>PATENT-JA-UPA-1995-000001</DOCNO>

<SDO BIJ>      <Bibliography>
(43)【公開日】平成７年（１９９５）……　　[Publication Date]
(54)【発明の名称】スラリ散布を行う……　　[Title of Invention]
……

<SDO ABJ>      <Abstract>
【目的】　スラリの処理と土壌作業を……　　[Purpose]

【構成】　トラクタとスラリを積載した……　　[Composition]

<SDO CLJ>      <Scope of Claim>
【請求項１】　バキウムカーを牽引して……　　[Claim1]
【請求項２】　トラクタに対して ……　　[Claim2]

<SDO DEJ>      <Description>
【産業上の利用分野】本発明はスラリ……　　[Application Field]

【発明が解決しようとする課題】このようなスラリを圃場に供給する……　　[Problem to be solved]

【課題を解決するための手段】上述のような目的を達成するために、……　　[Means of solving Problems]

【作用】本発明のスラリ散布を行う土壌作業機は、……　　[Operation]

【実施例】以下、本発明を採用した土壌作業機について添付した図面に……　　[Embodiment Example]

【発明の効果】以上の説明から明らかな……　　[Effects of Invention]

<SDO EDJ>      <Explanation of Drawings>
【図１】本発明のスラリ散布を行う……　　[Figure1]

【図１】……

<DOCNO>PATENT-JA-UPA-1995-000001</DOCNO>

<Technical Field>

<Purpose>

<Method>

<Claim>

<Explanation>

<Example>

These 6 fields are decided through the observation of applicant-defined tags

# System Architecture for Retrieval and Categorization

**Phase 1:** Document Indexing

Target patent documents
↓
Re-organization
↓
$Field_1$ $Field_2$ ... $Field_n$
↓
Keyword extraction
↓
Indexing
↓
$Index_1$ $Index_2$ ... $Index_n$

**comparison**

**Phase 2:** Document Retrieval

Query patent document
↓
Re-organization
↓
$Field_1$ $Field_2$ ... $Field_n$
↓
Keyword extraction
↓
$Query_1$ $Query_2$ ... $Query_n$
↓
Retrieval result$_1$ Retrieval result$_2$ ... Retrieval result$_n$
↓
Merging
↓
Document list

**Phase 3:** Document Categorization

Category list ← Category score calculation

# *How to re-organize a document?*

❖ Re-organizing by using applicant-defined tags

- Various tags in <ABJ>, <DEJ>
  - 3,516 tags (among 347,227 doc.)

| Frequency | applicant defined tags (Japanese) | Meaning |
|---|---|---|
| 310,276 | 課題を解決するための手段 | Means of solving problems |
| 2,502 | 問題点を解決するための手段 | |
| 1,449 | 課題を解決する為の手段 | |
| 2,923 | 課題を解決するための手段及び作用 | Means of solving problems & Operation |
| 3,962 | 従来の技術及び発明が解決しようとする課題 | Prior art & Problem to be solved |
| 306,350 | 発明が解決しようとする課題 | Problem to be solved |
| 2,121 | 発明が解決しようとする問題点 | |
| 1,476 | 発明が解決しようとしている課題 | |

- Classifying tags according to head nouns of tags
  - 100 most frequent HNs ➜ 6 classes

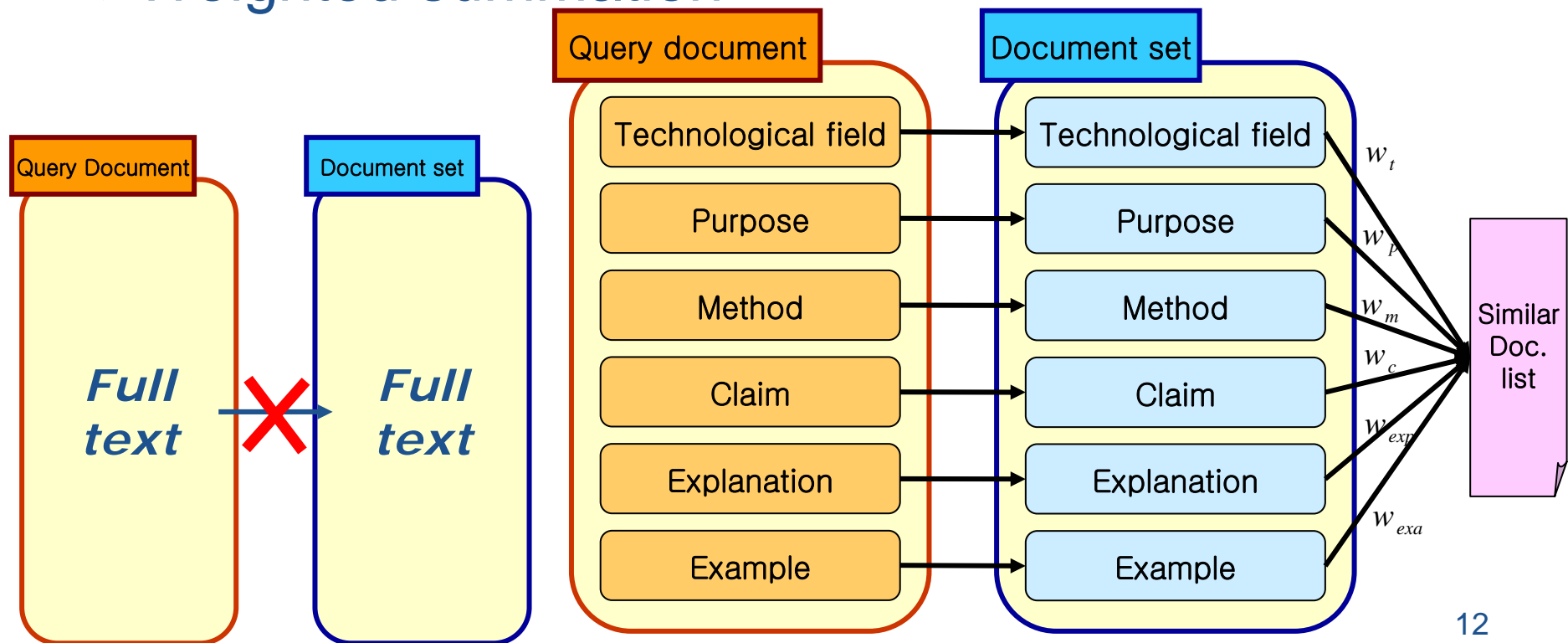10

# Examples of Classified Applicant-defined Tags

| Primitive label | Head | Examples of Applicant –defined tags |
|---|---|---|
| Technological field | 分野 (field)<br>技術 (art)<br>背景 (background) | 産業上の利用分野 (Industrial application field)<br>従来の技術 (prior art)<br>発明の背景 (background of the invention) |
| Purpose | 名称 (title)<br>目的 (purpose)<br>課題 (problem) | 発明の名称 (title of the invention)<br>発明の目的 (purpose of the invention)<br>発明が解決しようとする課題 (problem to be solved by the invention) |
| Method | 手段 (the means) | 問題点を解決するための手段 (the means of solving the problem)<br>課題を解決するための手段及び作用<br>(the means of solving the problem and the operation) |
| Claim | | All titles in the <Claim> part |
| Explanation | 構成 (composition)<br>効果 (effect)<br>作用 (operation)<br>説明 (explanation) | 構成 (Composition)<br>発明の効果 (the effect of the invention)<br>課題を解決するための手段及び作用<br>(the means of solving the problem and the operation)<br>発明の具体的説明 (The concrete explanation of composition) |
| Example | 例 (example) | 実施例 (embodiment example)<br>参考例 (referential example)<br>実験例 (experimental example) |

**Multiple classification**

11

**Phase 2: Retrieval**

❖ Pairs of same fields are compared

- 6 queries – 6 indexes
- by Lemur toolkit

❖ Weighted summation

| Query Document | Document set |
|---|---|
| **Full text** | **Full text** |

Query document

| Technological field | → | Technological field | $w_t$ |
| Purpose | → | Purpose | $w_p$ |
| Method | → | Method | $w_m$ |
| Claim | → | Claim | $w_c$ |
| Explanation | → | Explanation | $w_{exp}$ |
| Example | → | Example | $w_{exa}$ |

Document set

Similar Doc. list

12

❖ Method

- by theme codes of *k* documents similar to a query document

  - Retrieved documents have theme codes

    - Example) similarity result for a query document

| rank | doc ID | document similarity | given theme codes |
|------|--------|---------------------|-------------------|
| 1 | d04 | 371.773 | 2B062 |
| 2 | d01 | 371.009 | 2B062, 2B304 |
| 3 | d02 | 370.981 | 2B072 |
| 4 | d03 | 370.863 | 2B304, 3L045, 3L055 |
| 5 | d09 | 370.800 | 3L045 |
| …… | | | |

*K=3 means that top 3 documents are meaningful among N retrieved documents*

13

# Assigning theme codes (2/2)

❖ Method for calculating scores of theme codes

**Example for a given query**

| Similarity Result | Score for Theme code |
|---|---|

**Weight value α = 1**

| doc rank | doc ID | document similarity | given Theme codes |
|---|---|---|---|
| 1 | d04 | 371.773 | 2B062 |
| 2 | d01 | 371.009 | 2B062  2B304 |
| 3 | d02 | 370.981 | 2B072 |
| 4 | d03 | 370.863 | 2B304, 3L045, 3L055 |
| 5 | d09 | 370.800 | 3L045 |
| …… | | | |

**k=3**

| theme rank | Theme code | score for theme code |
|---|---|---|
| 1 | 2B062 | 371.773+371.009 |
| 2 | 2B304 | 371.009+370.863*0.1 |
| 3 | 2B072 | 370.981 |
| 4 | 3L045 | 370.863*0.1+370.8*0.1 |
| 5 | 3L055 | |
| …… | | |

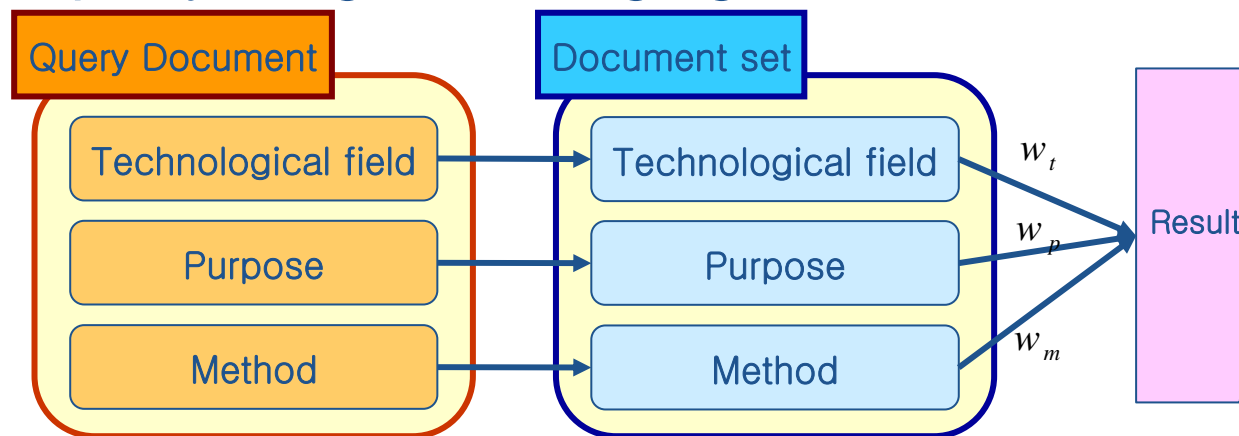**Weight value α = 0.1** until N(=200)th rank

14

# Experimental Results
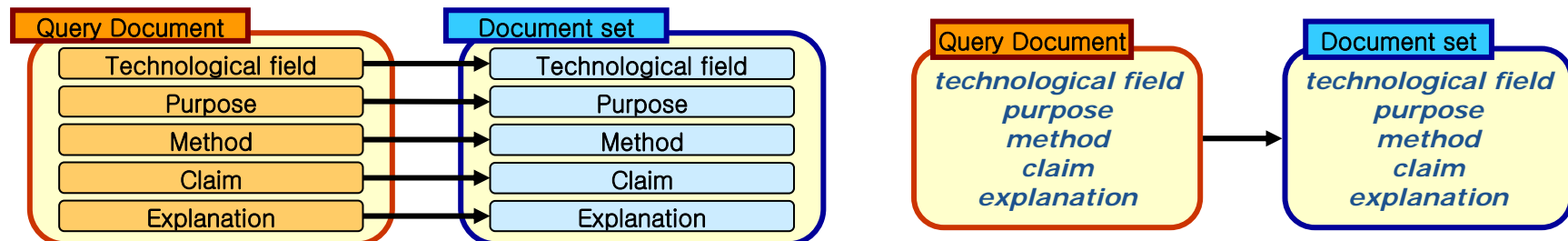
# *Theme Categorization Subtask*

❖ kNN-based Theme categorization

- Training documents: 2 years documents (1993, 1997)
- Component-by-component retrieval
- Equally weighted merging

| Query Document | Document set | | Result |
|---|---|---|---|
| Technological field | Technological field | $w_t$ | |
| Purpose | Purpose | $w_p$ | |
| Method | Method | $w_m$ | |

| RunID | Condition | MAP | Rank |
|---|---|---|---|
| ft001 | k=10 | 0.6872 | 1 |
| ft002 | k=20 | 0.6842 | 2 |
| ft003 | k=30 | 0.6819 | 3 |

# *Additional experiments*

❖ kNN-based approach vs. Word-bag Vector

- ▪ MAP in NTCIR-5 formal run
  - 0.6872 (kNN) > 0.3776 (MEM)

❖ Detailed component vs. Normative section

- ▪ Dev. Set (1 year training data, 100 test documents)
  - 0.6372 (Purpose) > 0.5774 (Description section)

❖ Component-by-Component vs. Doc–by-Doc

- ▪ Dev. Set
  - 0.6402 (Technical field, Purpose, Method, Claim, Explanation)
  - 0.6050 (Technical field+Purpose+Method+Claim+Explanation)

| Query Document | | Document set |
| --- | --- | --- |
| Technological field | → | Technological field |
| Purpose | → | Purpose |
| Method | → | Method |
| Claim | → | Claim |
| Explanation | → | Explanation |

| Query Document | Document set |
| --- | --- |
| *technological field* *purpose* *method* *claim* *explanation* | *technological field* *purpose* *method* *claim* *explanation* |

# *Conclusions*

❖ kNN based approach for patent categorization

- Using similar N documents
  - ➔ the effect of feature reduction (cf. word-bag vector)

❖ Component-by-component comparison

- Considering meaningful small units in a document
  - ➔ precise comparison of content among documents
  - cf. document-by-document comparison