

How to Effectively Use Topic Models for Software Engineering Tasks? An Approach Based on Genetic Algorithms



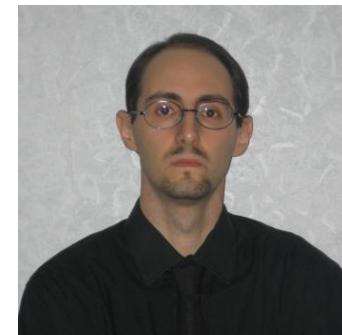
Annibale
Panichella



Bogdan
Dit



Rocco
Oliveto



Massimiliano
Di Penta



Denys
Poshyvanyk



Andrea
De Lucia





Source Code

```

Private Function CleanUpLine(ByVal sLine As String) As String
Dim iLocnCCount As Long
Dim iLocnLCount As Long
Dim sChar As String
Dim sPrevChar As String

' Starts with Rem it is a comment
sLine = Trim(sLine)
If Left(sLine, 3) = "Rem" Then
    CleanUpLine = ""
    Exit Function
End If

' Starts with ' it is a comment
If Left(sLine, 1) = "'" Then
    CleanUpLine = ""
    Exit Function
End If

' Contains '' may end in a comment, so test if it is a comment or in the
' body of a string
If InStr(sLine, "") > 0 Then
    sPrevChar = ""
    iQuoteCount = 0

    For lCount = 1 To Len(sLine)
        sChar = Mid(sLine, lCount, 1)

        ' If we found "" then an even number of " characters in front
        ' means it is the start of a comment, and odd number means it is
        ' part of a string
        If iLocnCCount Mod 2 = 0 Then
            If sChar = "" Then
                sLine = Trim(Left(sLine, lCount - 1))
                Exit For
            End If
        ElseIf sChar = "" Then
            iQuoteCount = iQuoteCount + 1
        End If
        sPrevChar = sChar
    Next lCount
    CleanUpLine = sLine
End Function

```

Bug Reports

First Last Prev Next | This bug is not in your last search results.

Bug 816298 - Change "-moz-user-select:none" to behave like WebKit, IE, and Opera (and "-moz-user-select:moz-none")

Status: RESOLVED FIXED
Reported: 2012-11-28 15:03 PST by Chris Peters
Modified: 2013-05-20 00:14 PDT (History)
CC List: 12 users (Show)

Keywords: compat, dev-doc-complete
Product: Core (show info)
Component: DOM: CSS Object Model (Show info)
Version: Trunk
Platform: All All
Importance: normal (vote)
Target Milestone: mozilla21
Assigned To: Chris Peterson (cpeterson)
QA Contact:
URL: https://developer.mozilla.org /en-US/d...

Depends on: 739398 828534
Blocks: 790029 814974
Show dependency tree
/graph

Documentation

Java API
Wakes up all threads that are waiting on this object's monitor.
String toString()
Returns a string representation of the object.
void wait()
Causes current thread to wait until another thread invokes the `notify()` method or the `notifyAll()` method for this object.
void wait (long timeout)
Causes current thread to wait until either another thread invokes the `notify()` method or the `notifyAll()` method for this object, or a specified amount of time has elapsed.
void wait (long timeout, int nanosecond)
Causes current thread to wait until another thread invokes the `notify()` method or the `notifyAll()` method for this object, or some other thread interrupt the current thread, or a certain amount of real time has elapsed.

Constructor Detail

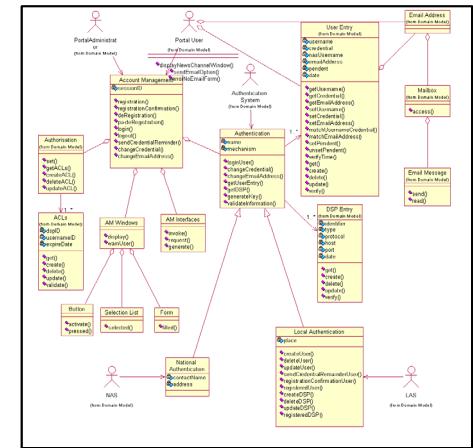
Object

Method Detail

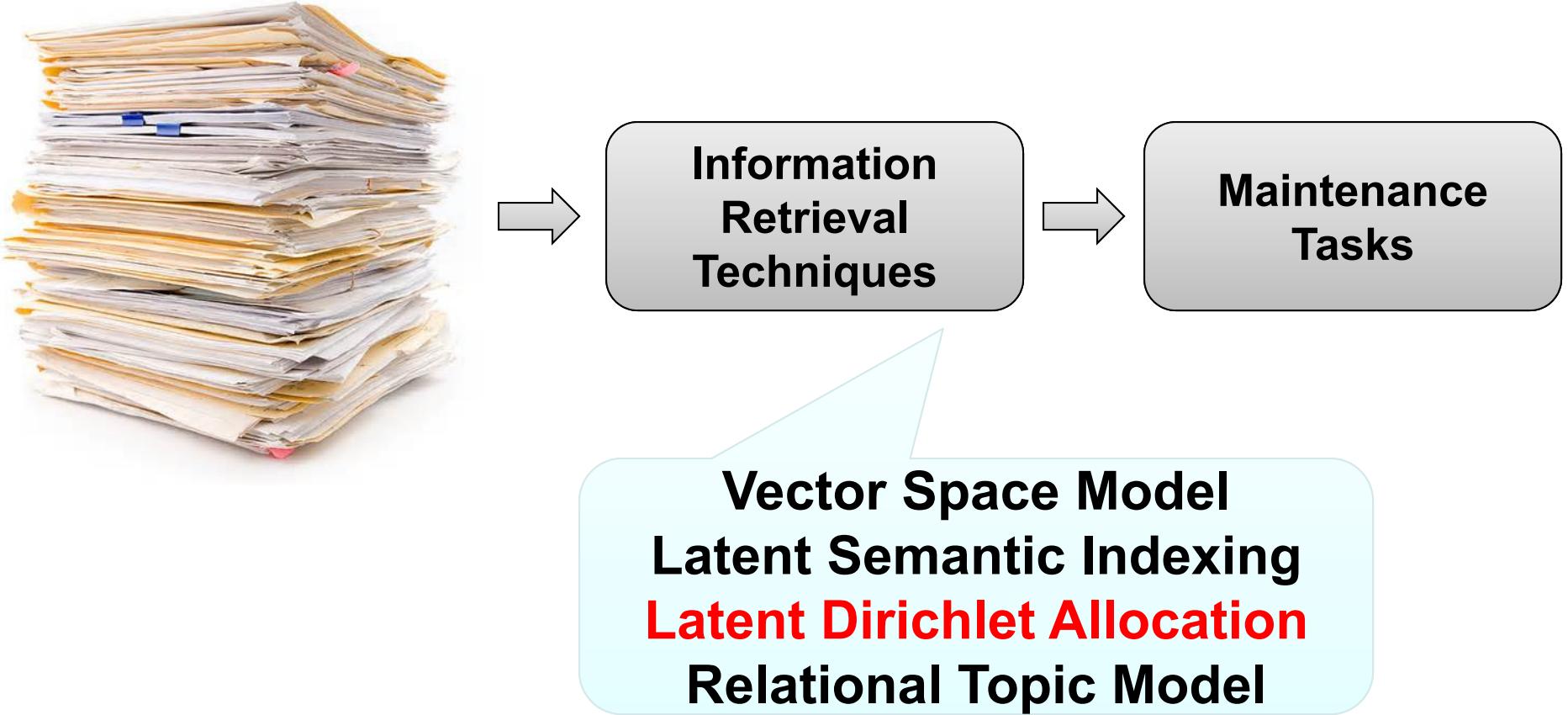
getClassName
Returns the runtime class of an object. That class is the object that is held by static synchronized methods of the represented class.

Returns:
The object of type Class that represents the runtime class of the object

Design documents







What is LDA?

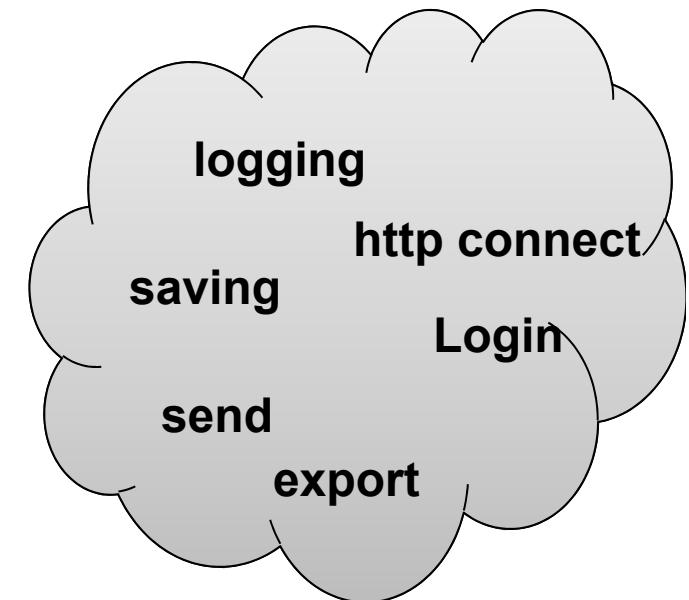
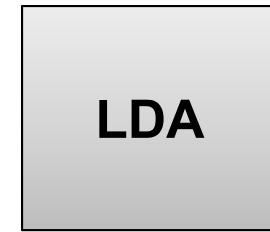
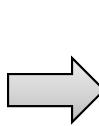
Latent Dirichlet Allocation (LDA)

- Topic model that generates the distribution of latent topics from textual documents

Latent Dirichlet Allocation (LDA)

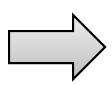
- Topic model that generates the distribution of latent topics from textual documents

```
Private Function CleanUpLine(ByVal sLine As String) As String
Dim lQuoteCount As Long
Dim sChar As String
Dim sP As String
    Private Function CleanUpLine(ByVal sLine As String) As String
Dim lQuoteCount As Long
Dim lCount As Long
Dim sChar As String
    Dim sPrevChar As String
    ' Starts with Rem it is a comment
    sLine = Trim(sLine)
    -If lCount > 0 Then
        CleanUpLine = sLine
    End If
    -Else
        -Starts with '
        sLine = Trim()
    End If
    Private Function CleanUpLine(ByVal sLine As String) As String
Dim lQuoteCount As Long
Dim lCount As Long
Dim sChar As String
Dim sPrevChar As String
    ' Starts with Rem it is a comment
    sLine = Trim(sLine)
    If Left(sLine, 3) = "Rem" Then
        CleanUpLine = ""
        Exit Function
    End If
    ' Starts with '
    If Left(sLine, 1) = "'" Then
        CleanUpLine = ""
        Exit Function
    End If
    ' Starts with it is a comment
    If Left(sLine, 1) = "*** Then
        CleanUpLine = ""
        Exit Function
    End If
    ' Contains ' means it is part of a string
    If InStr(sLine, "'") > 0 Then
        sChar = Mid(sLine, 1)
        lQuoteCount = 0
        -For lCount = 1 To Len(sLine)
            sChar = Mid(sLine, lCount, 1)
            -If sChar = "'" Then
                sLine = Left(sLine, lCount - 1)
                Exit For
            -ElseIf sChar = "" And sPrevChar = " "
                lQuoteCount = lQuoteCount + 1
            -End If
            sPrevChar = sChar
        -Next lCount
        CleanUpLine = sLine
    End If
    CleanUpLine = sLine
End Function
```



Software Artifacts

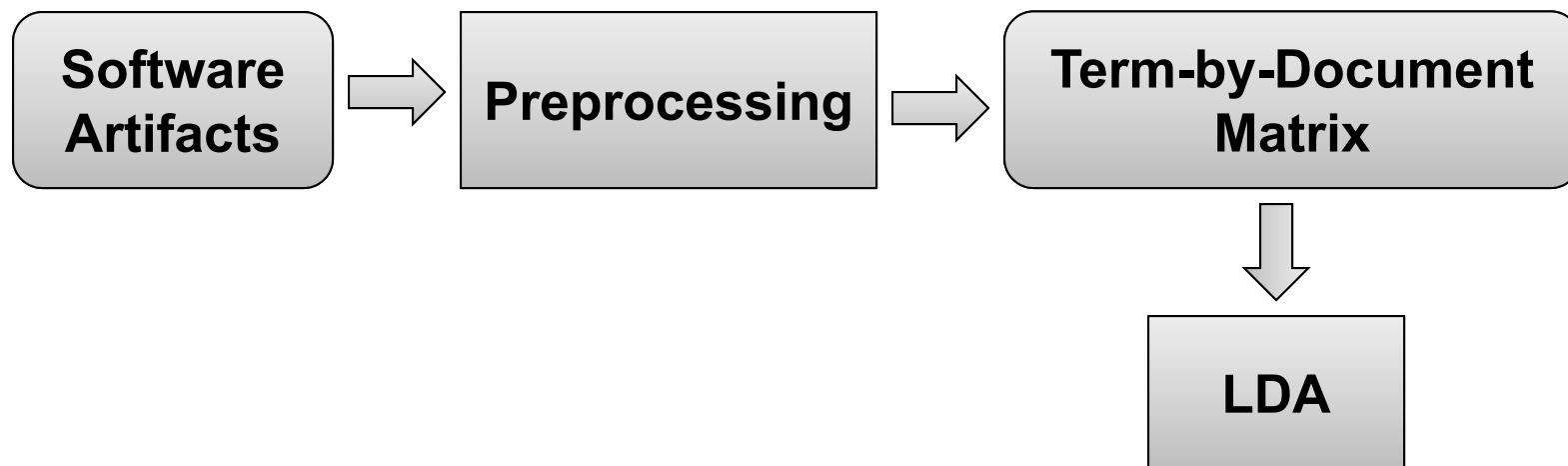
**Software
Artifacts**

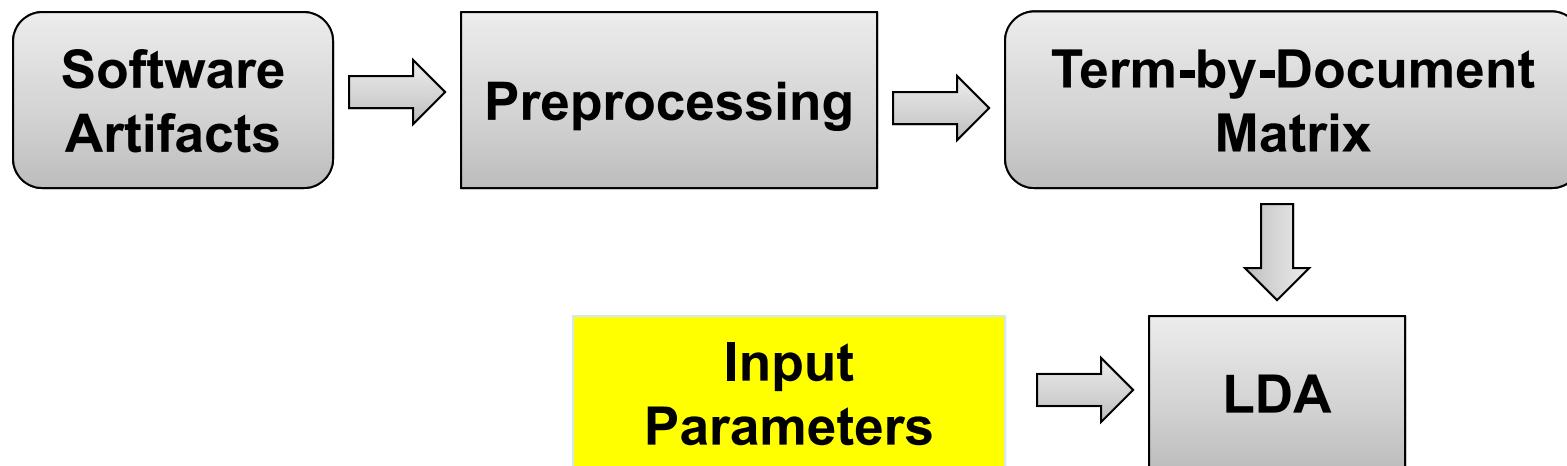


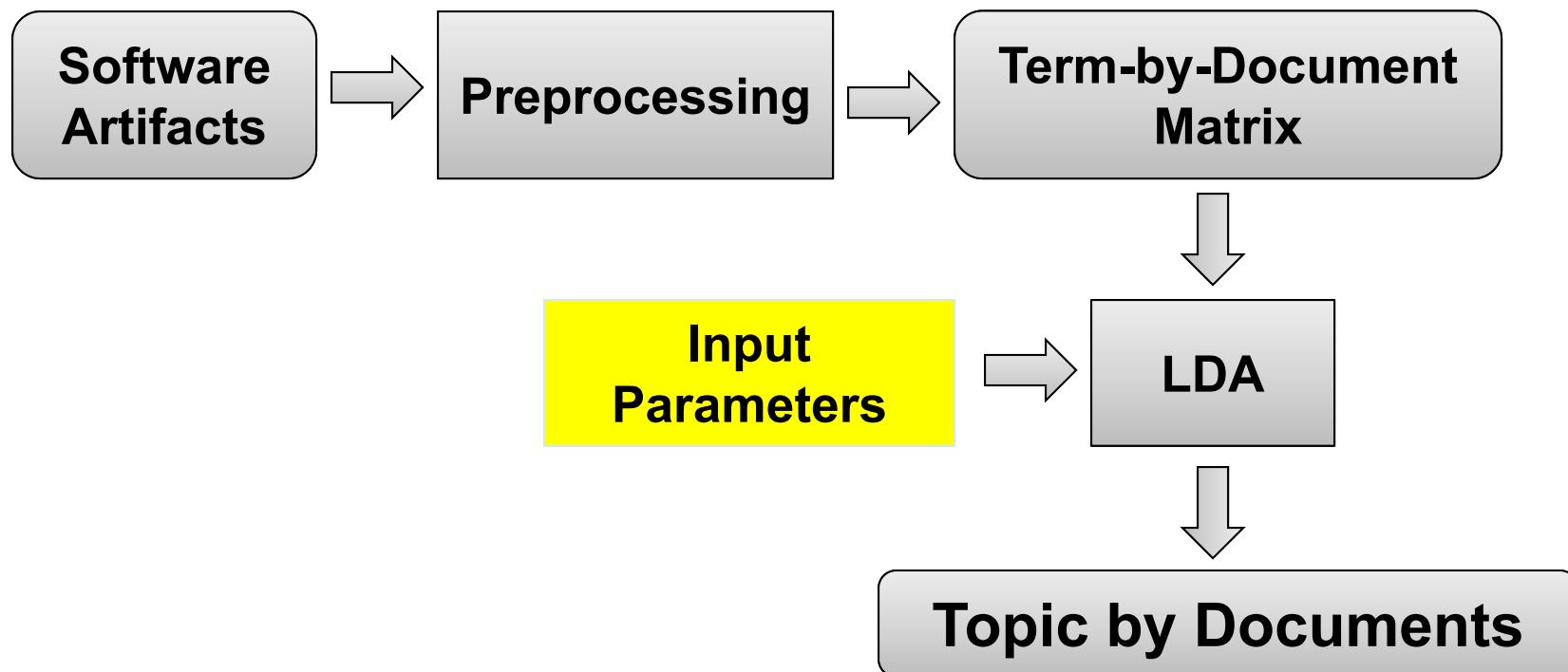
Preprocessing

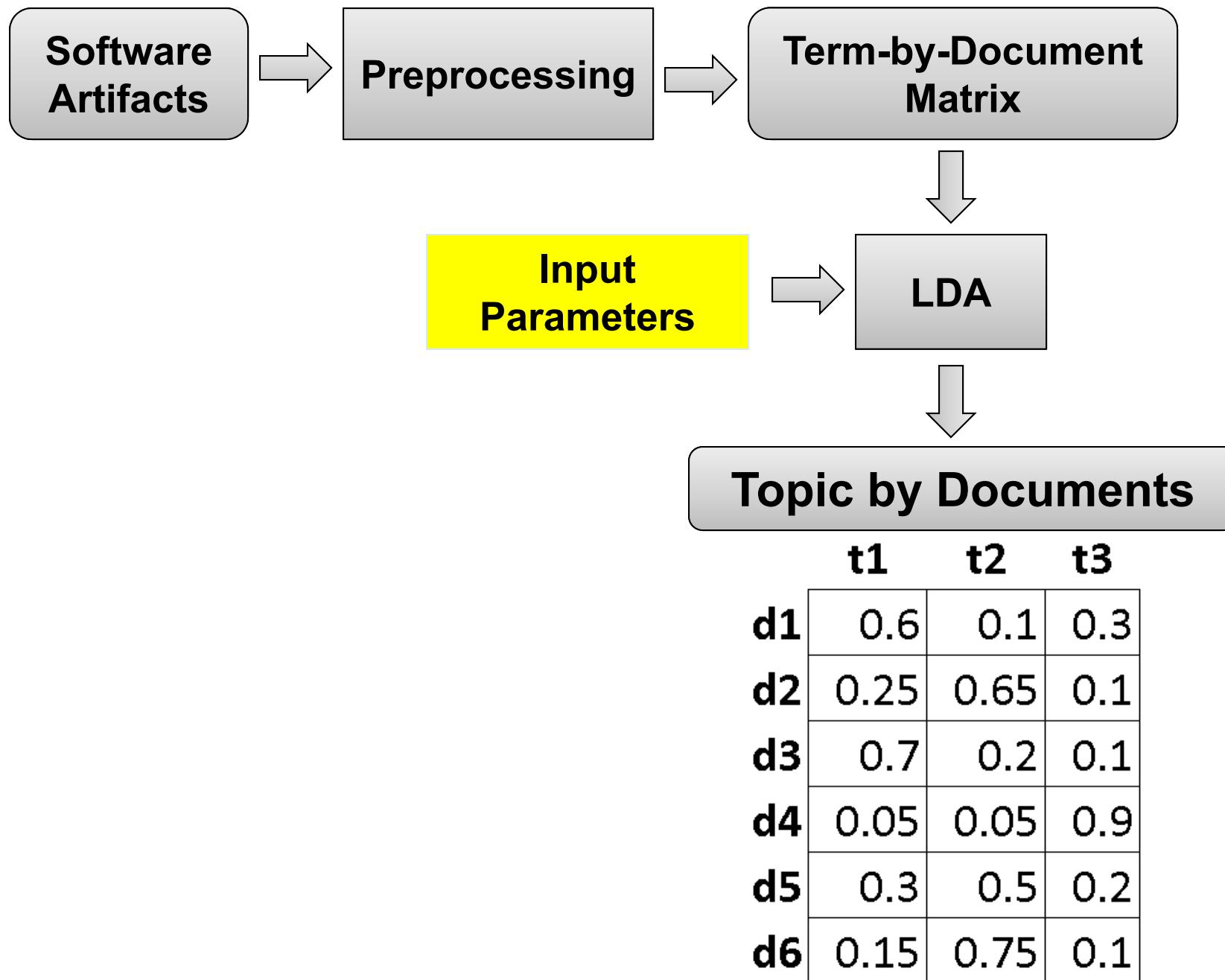
Remove special characters
Split identifiers
Remove common words
Stem

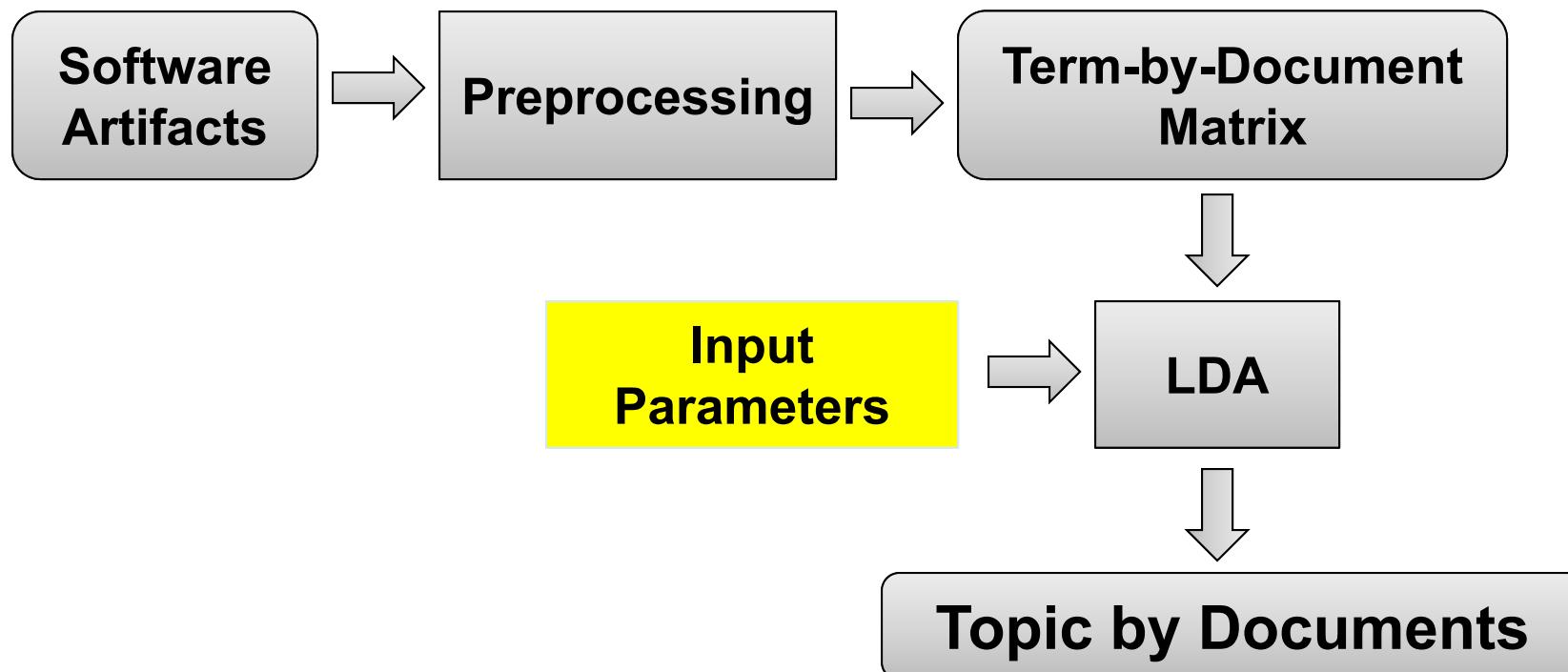






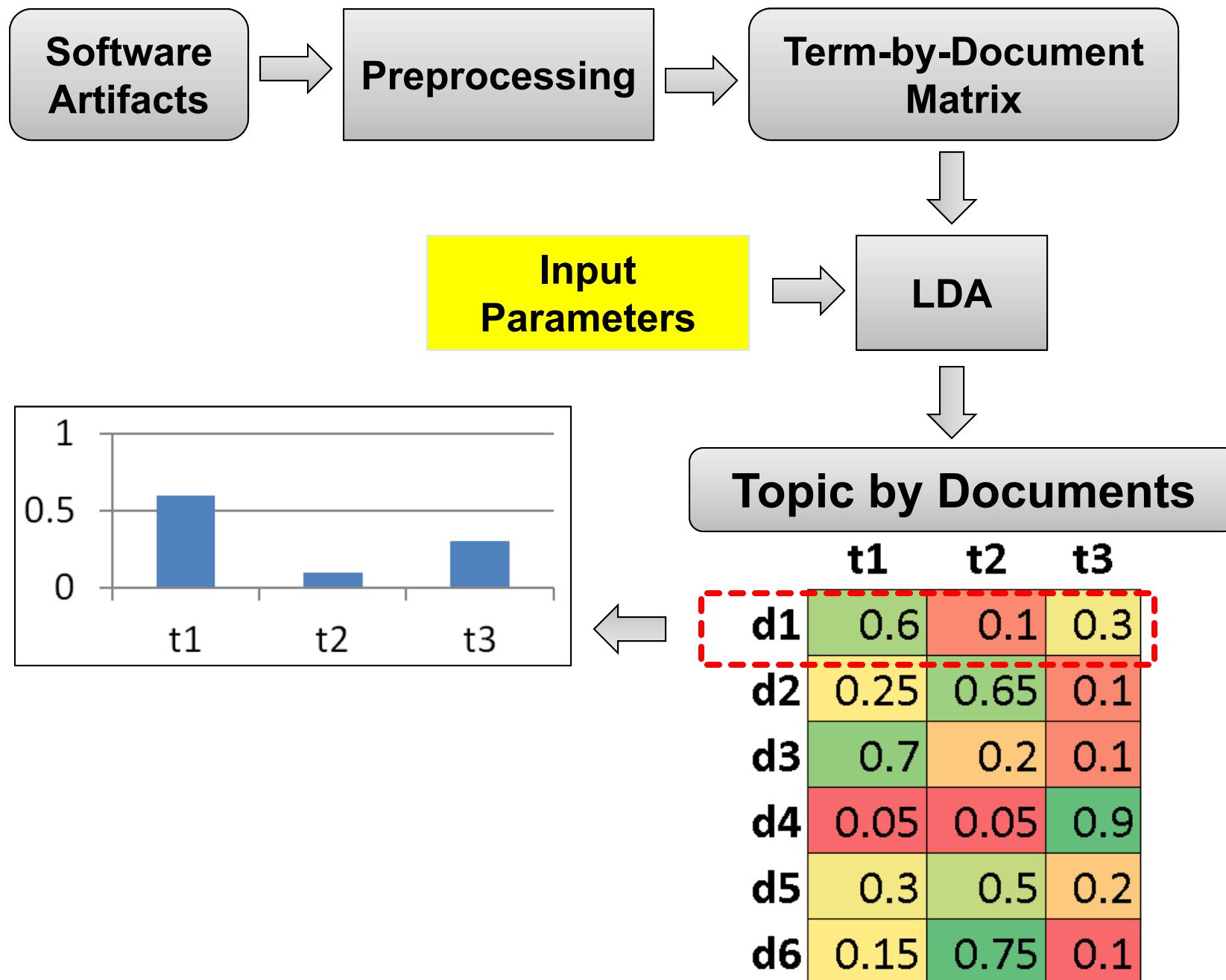


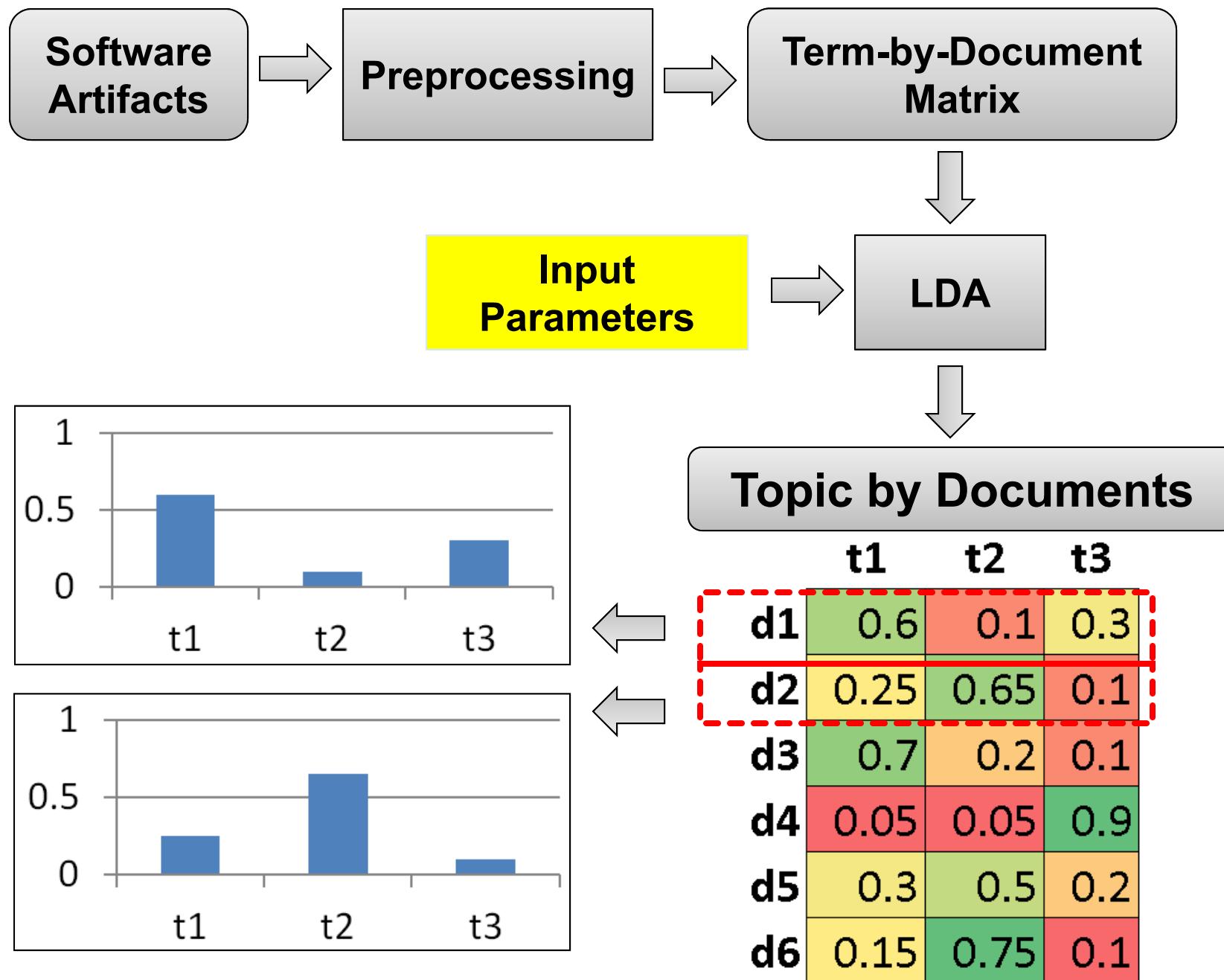


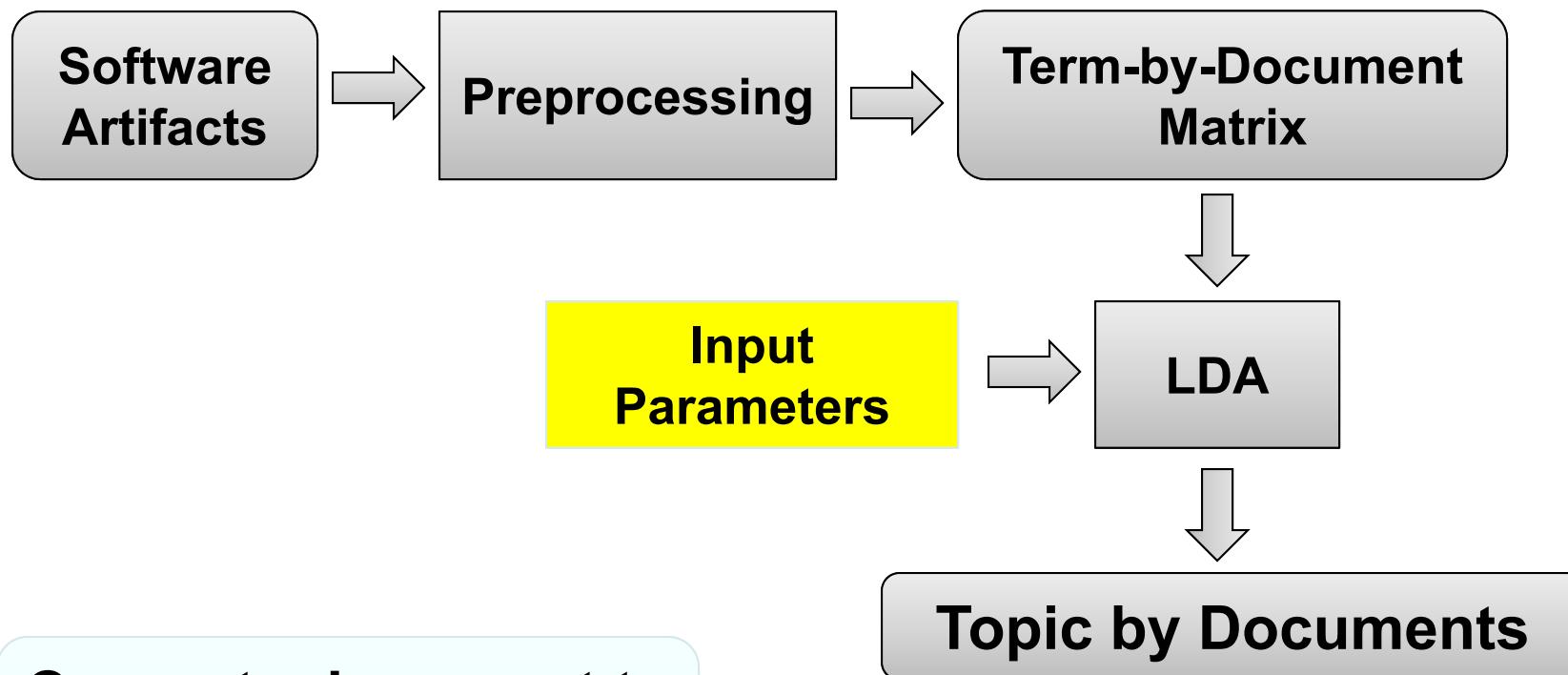


Probability that document is related to topic

| | t1 | t2 | t3 |
|----|------|------|-----|
| d1 | 0.6 | 0.1 | 0.3 |
| d2 | 0.25 | 0.65 | 0.1 |
| d3 | 0.7 | 0.2 | 0.1 |
| d4 | 0.05 | 0.05 | 0.9 |
| d5 | 0.3 | 0.5 | 0.2 |
| d6 | 0.15 | 0.75 | 0.1 |







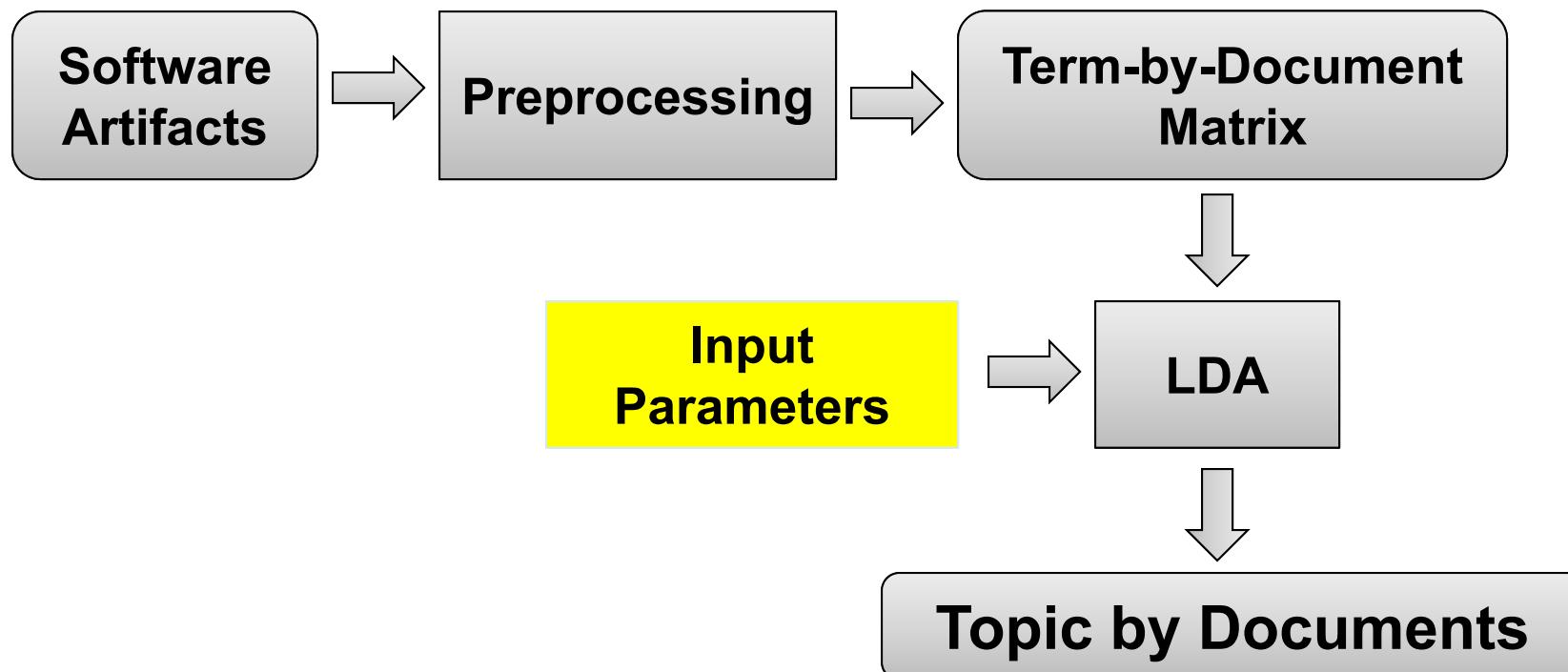
Compute document to document similarity

Compute query to document similarity

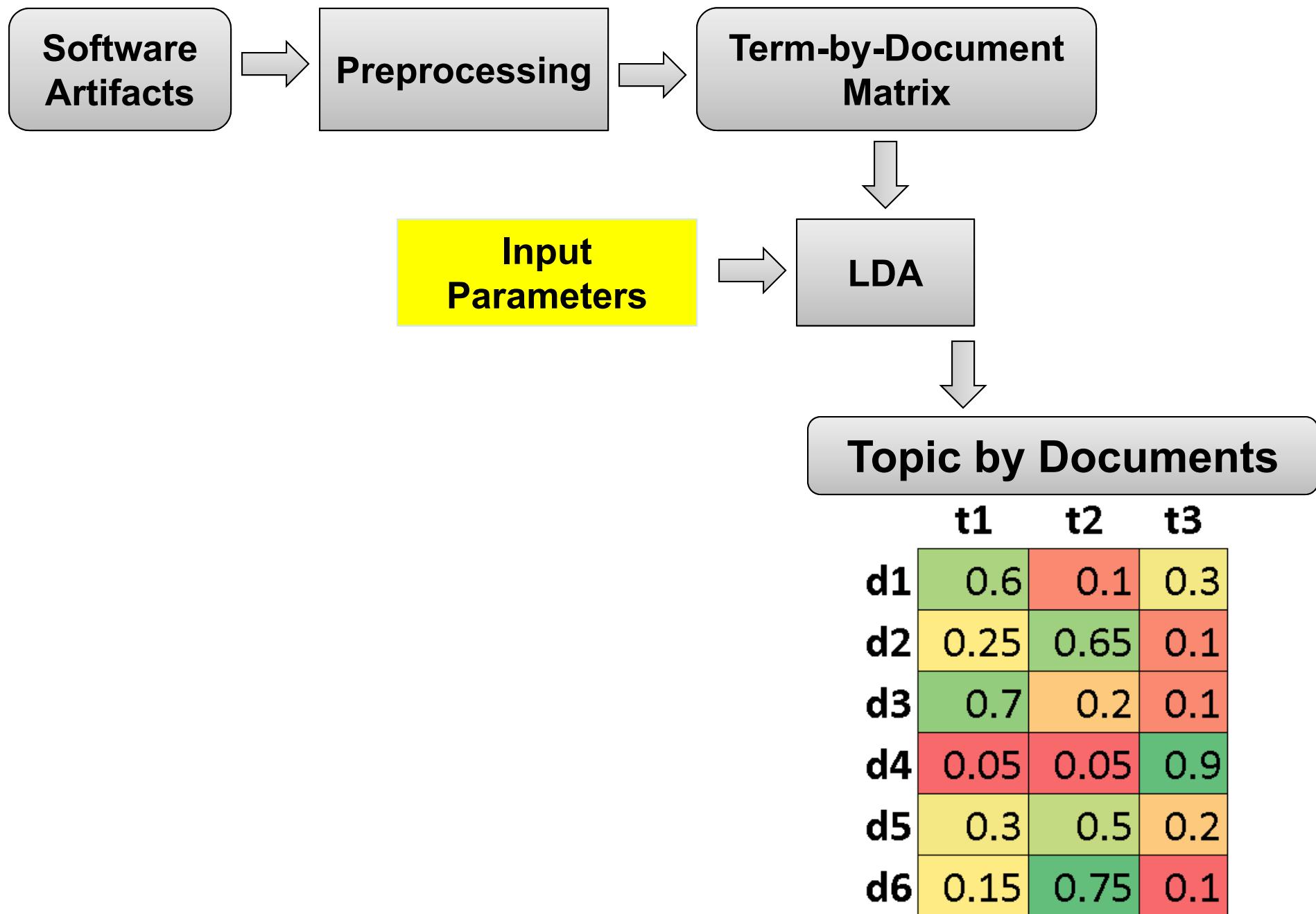
“Cluster” documents by topics

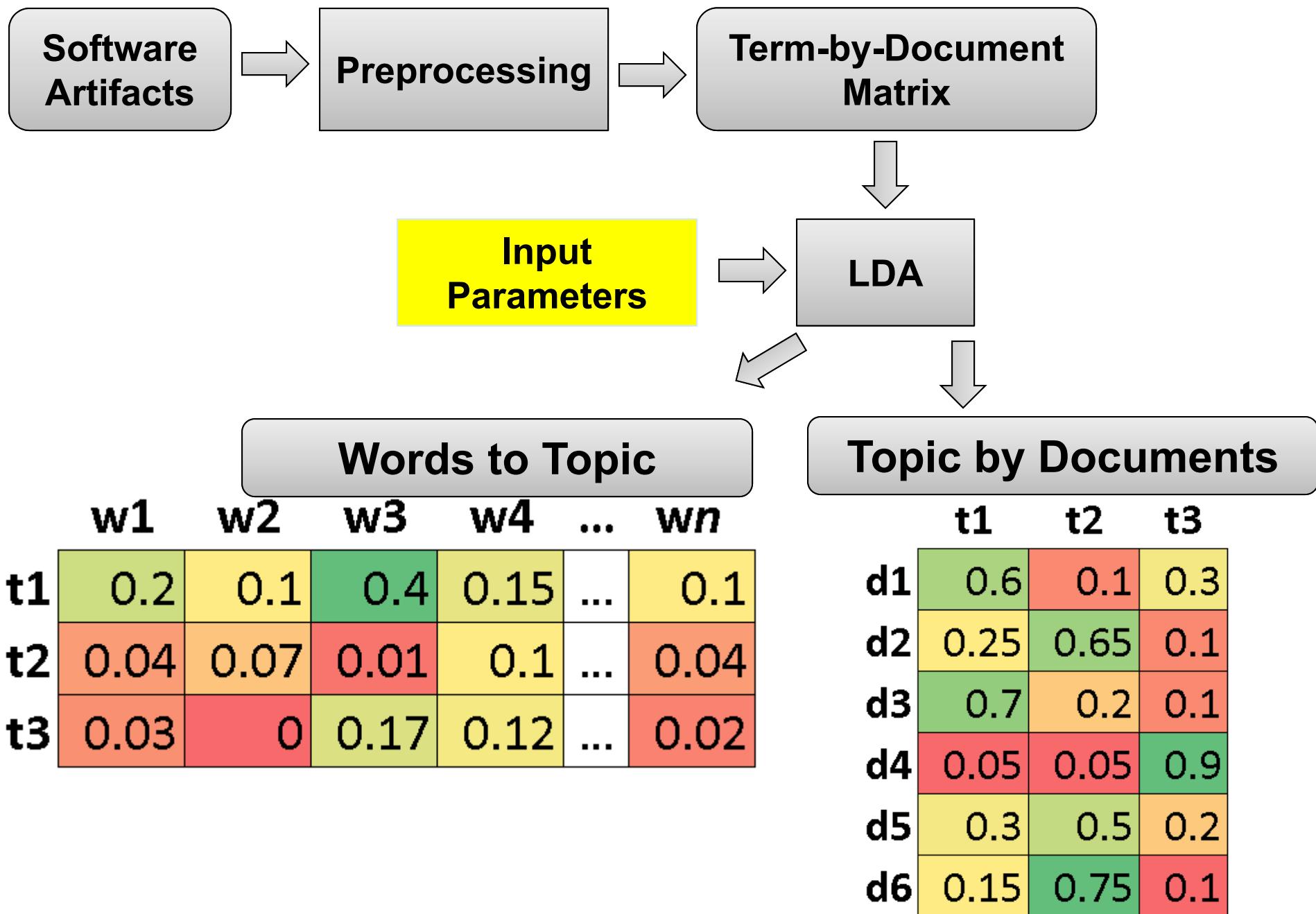
Topic by Documents

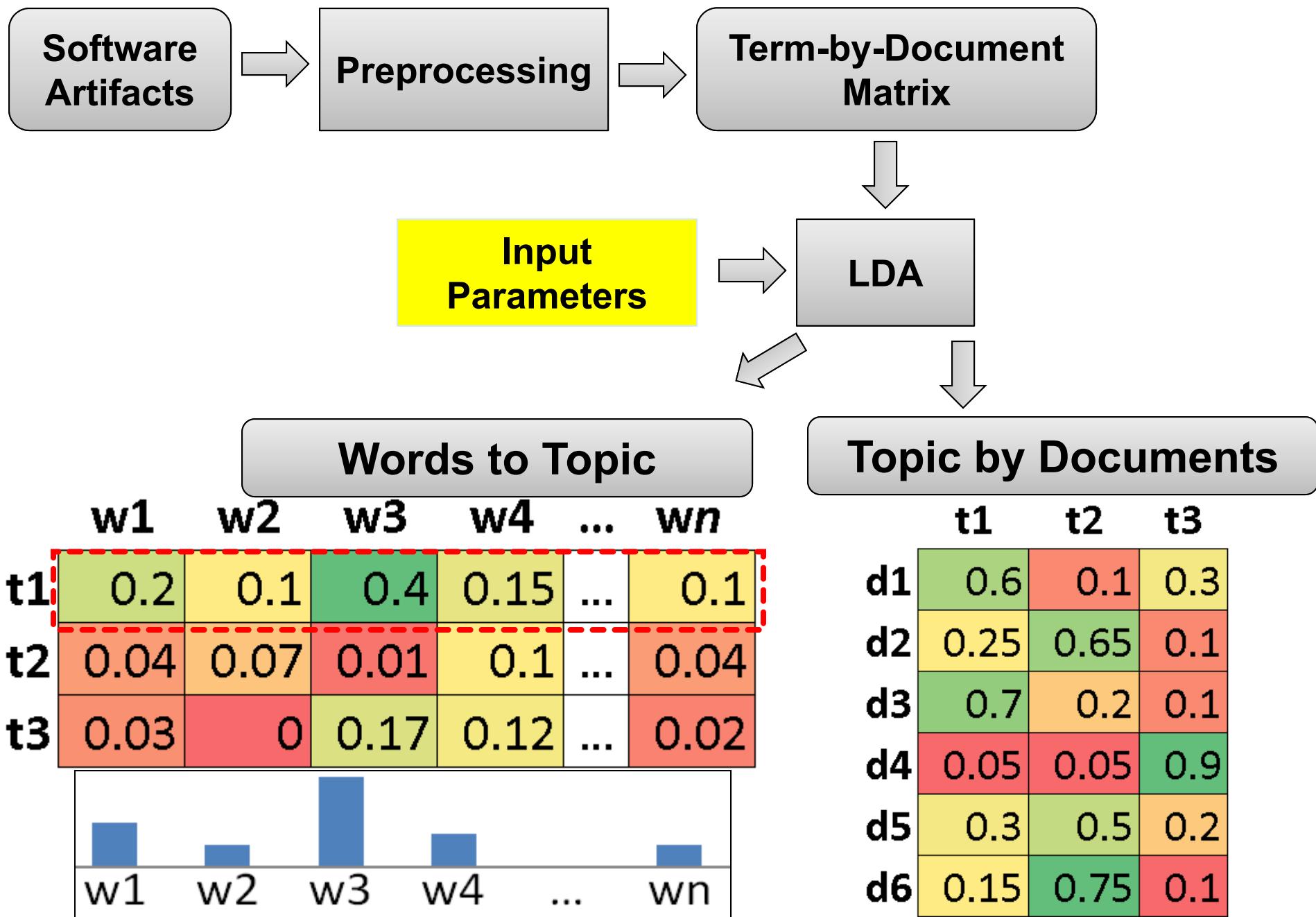
| | t1 | t2 | t3 |
|----|------|------|-----|
| d1 | 0.6 | 0.1 | 0.3 |
| d2 | 0.25 | 0.65 | 0.1 |
| d3 | 0.7 | 0.2 | 0.1 |
| d4 | 0.05 | 0.05 | 0.9 |
| d5 | 0.3 | 0.5 | 0.2 |
| d6 | 0.15 | 0.75 | 0.1 |

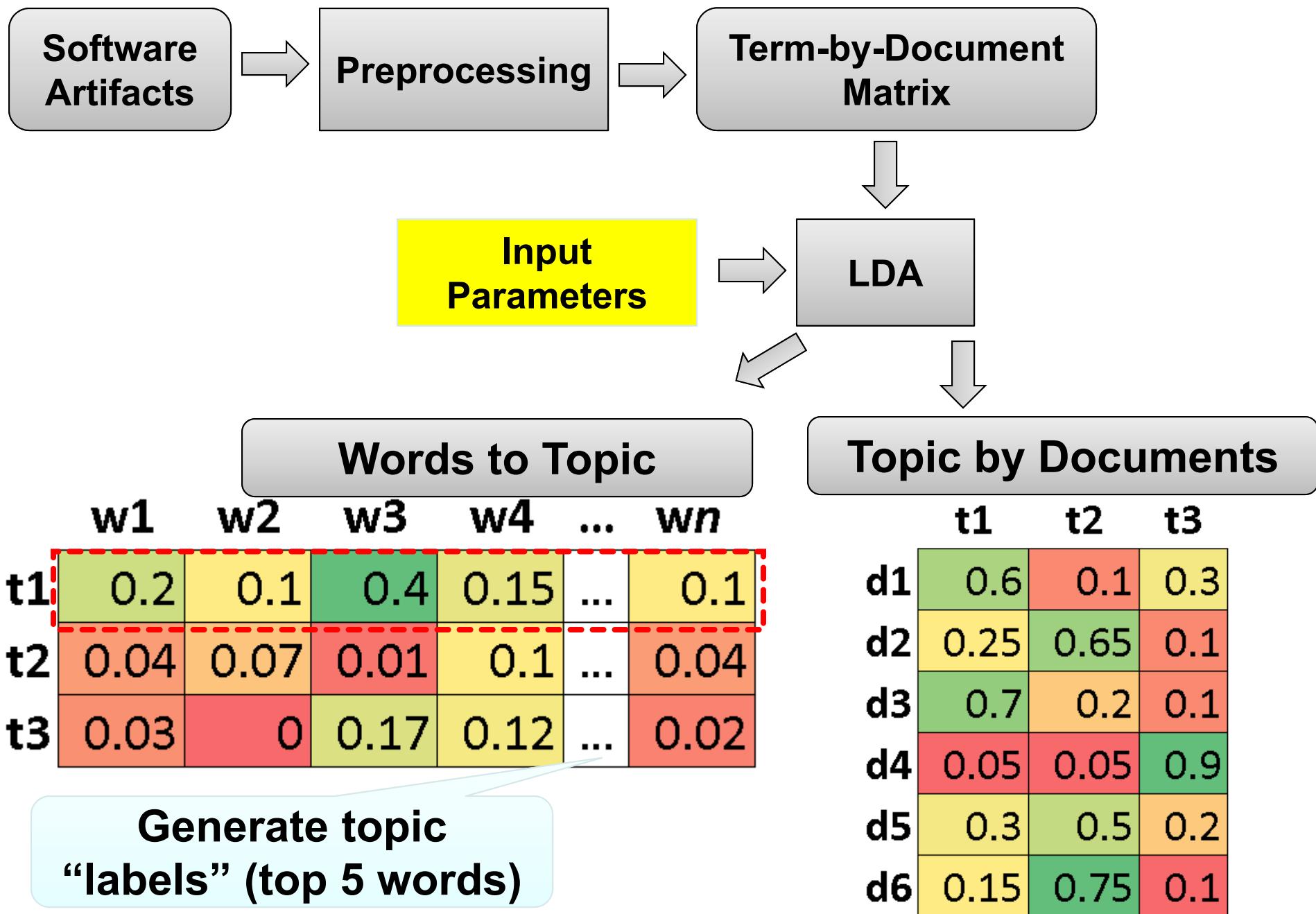


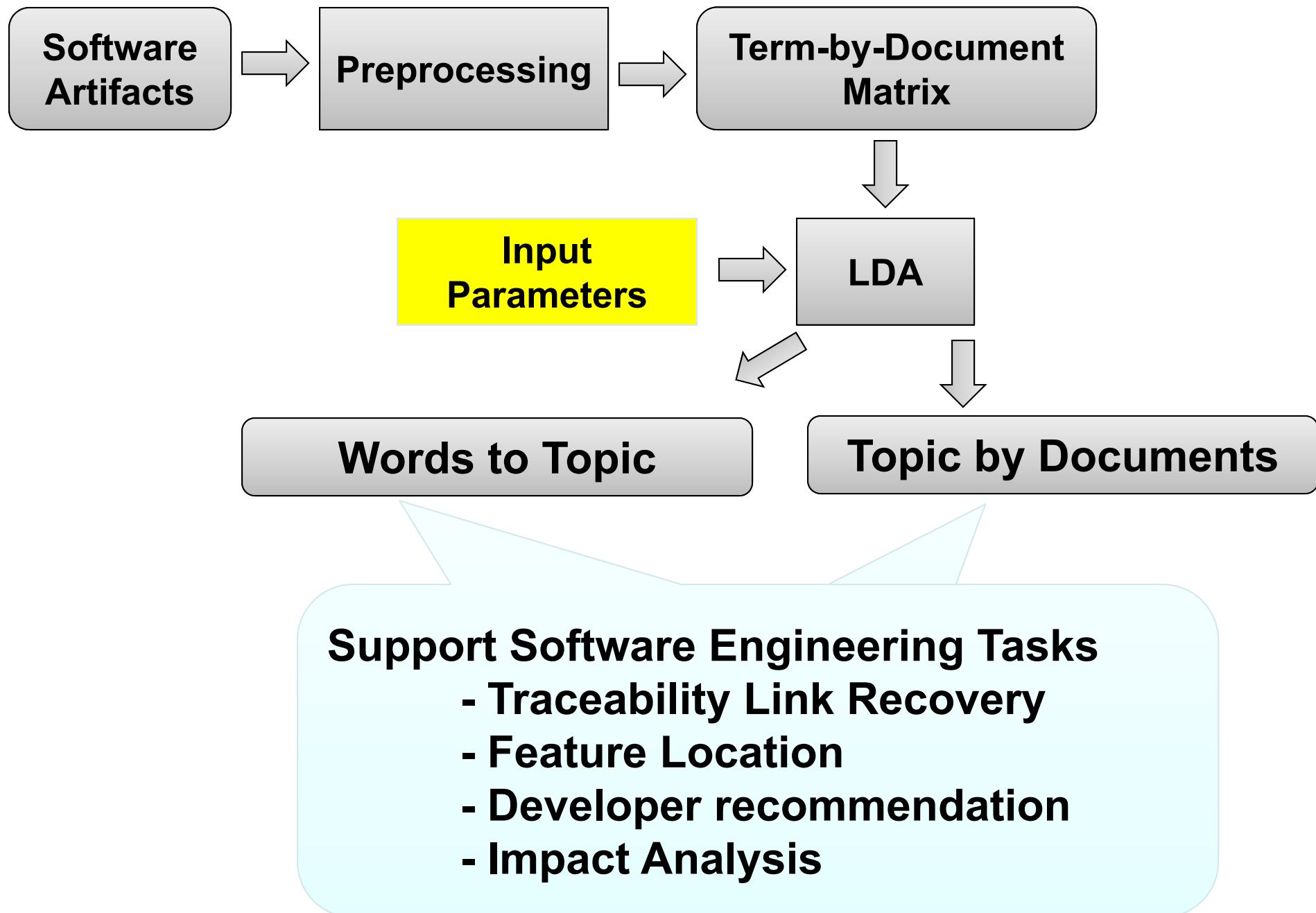
| | t1 | t2 | t3 |
|----|------|------|-----|
| d1 | 0.6 | 0.1 | 0.3 |
| d2 | 0.25 | 0.65 | 0.1 |
| d3 | 0.7 | 0.2 | 0.1 |
| d4 | 0.05 | 0.05 | 0.9 |
| d5 | 0.3 | 0.5 | 0.2 |
| d6 | 0.15 | 0.75 | 0.1 |

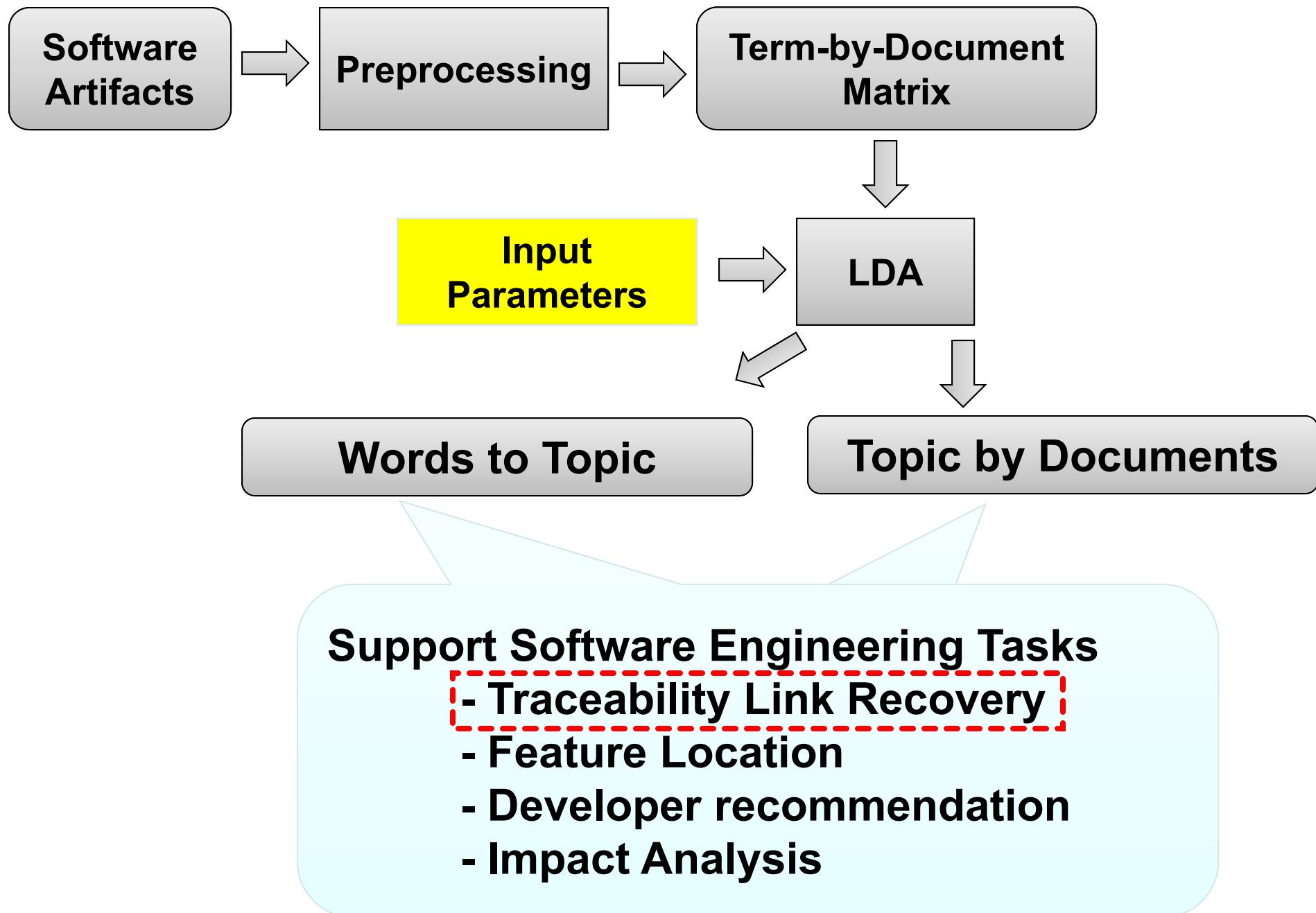








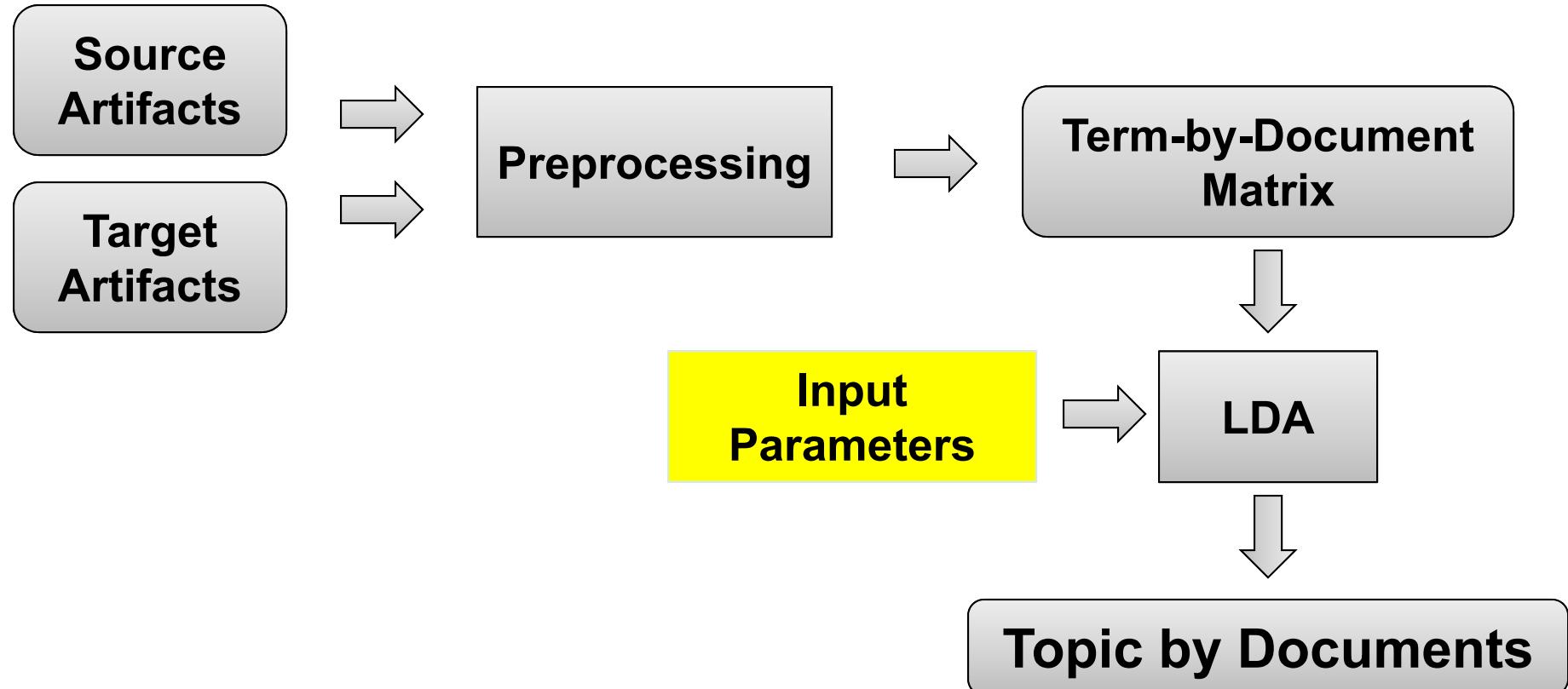


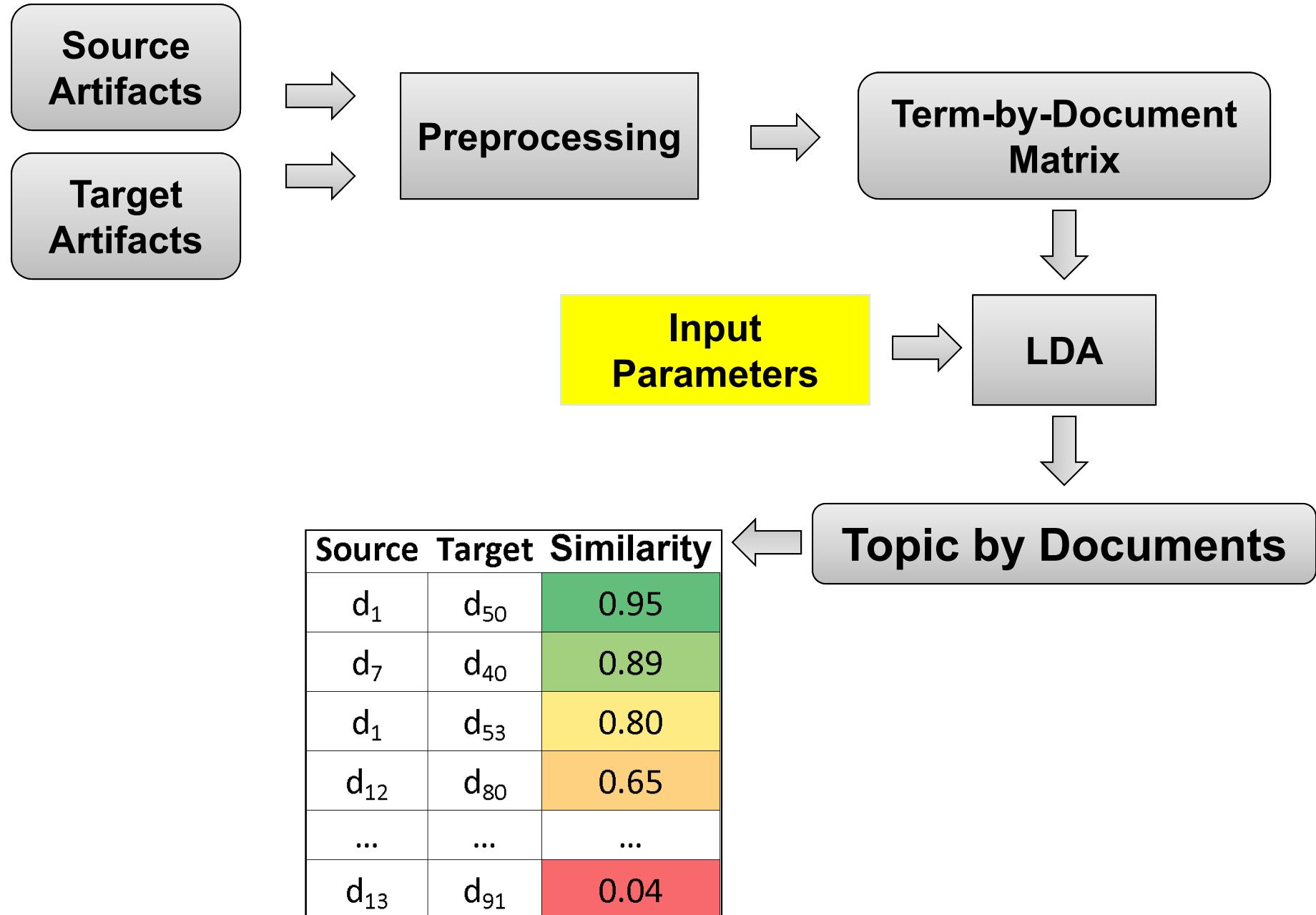


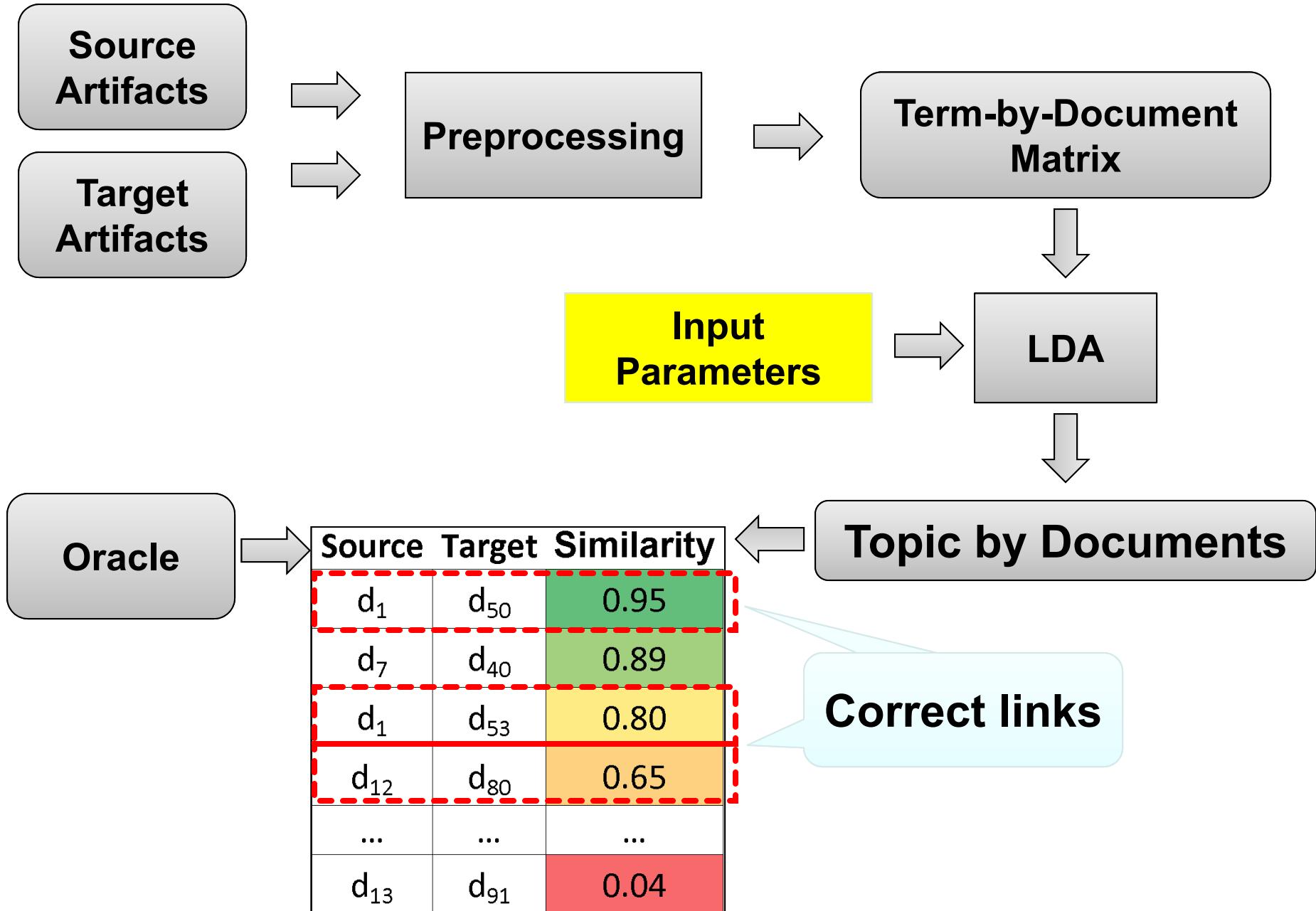
**Source
Artifacts**

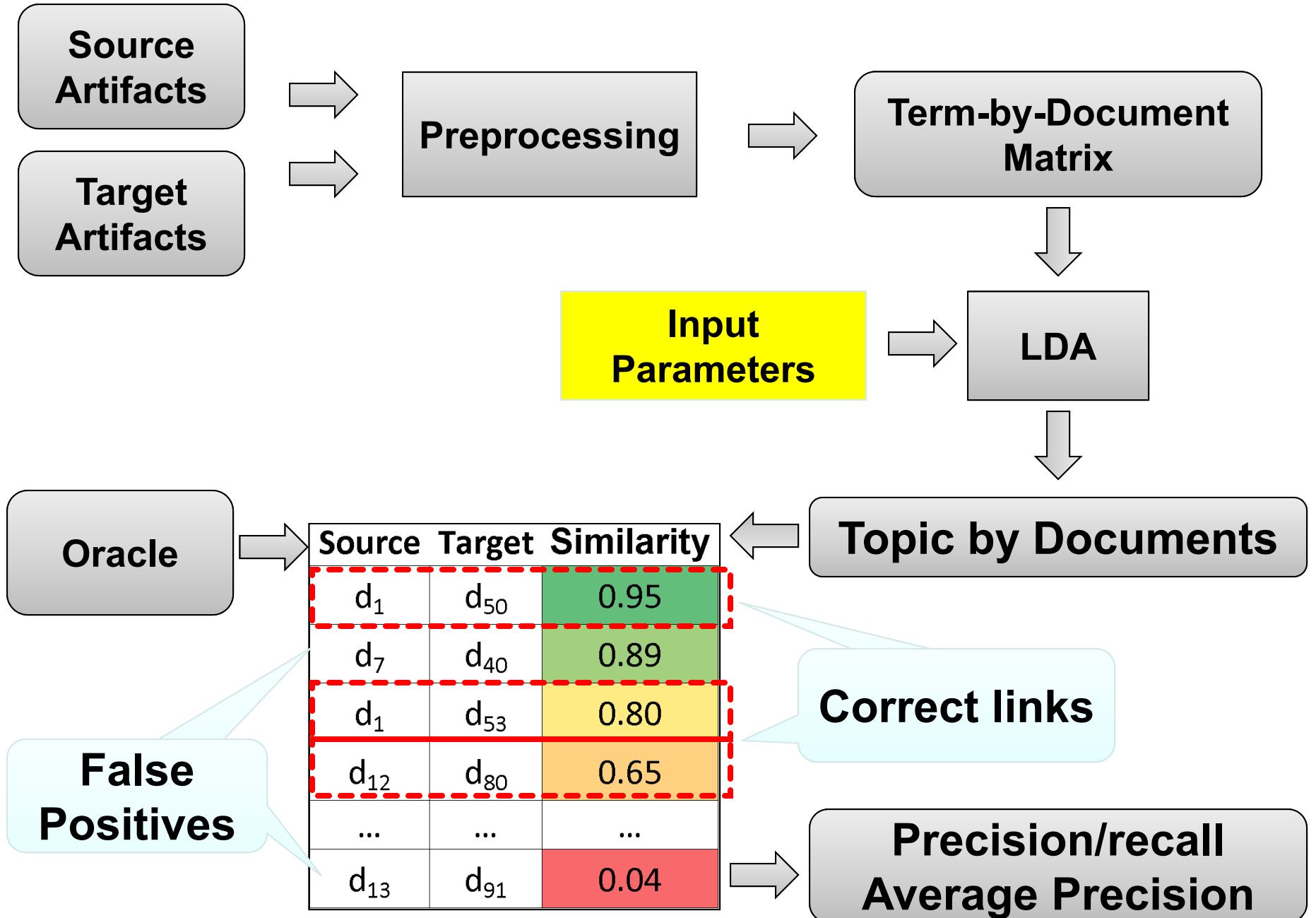
**Target
Artifacts**

**Use cases,
requirements,
classes, etc.**

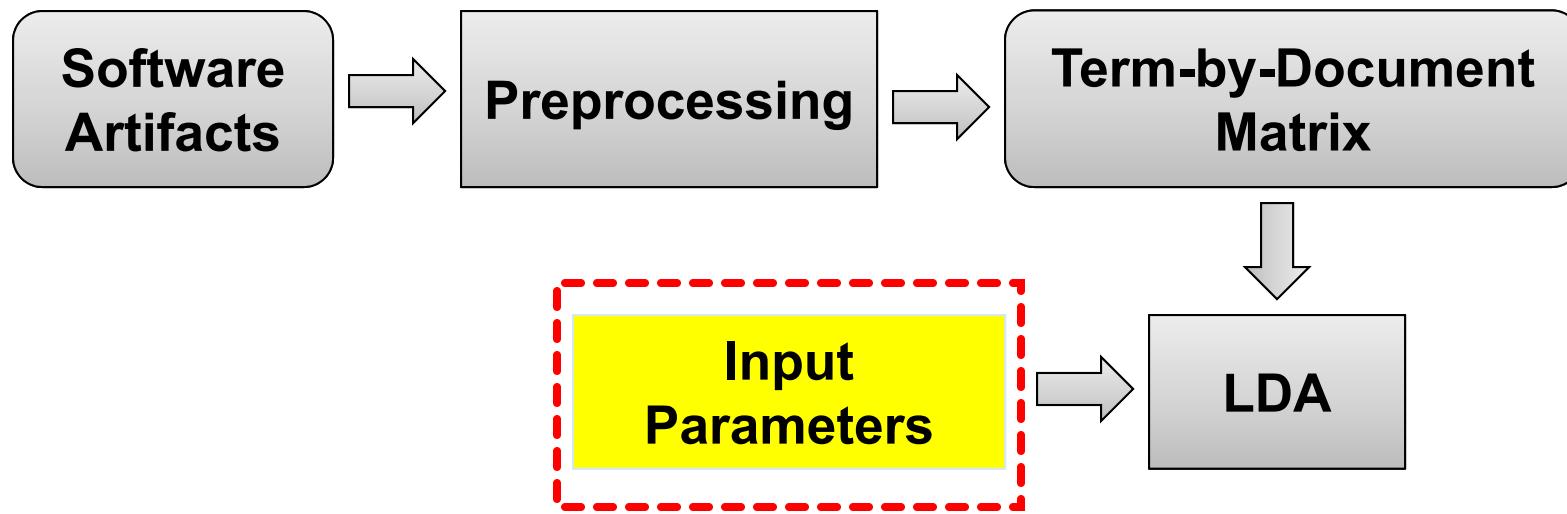


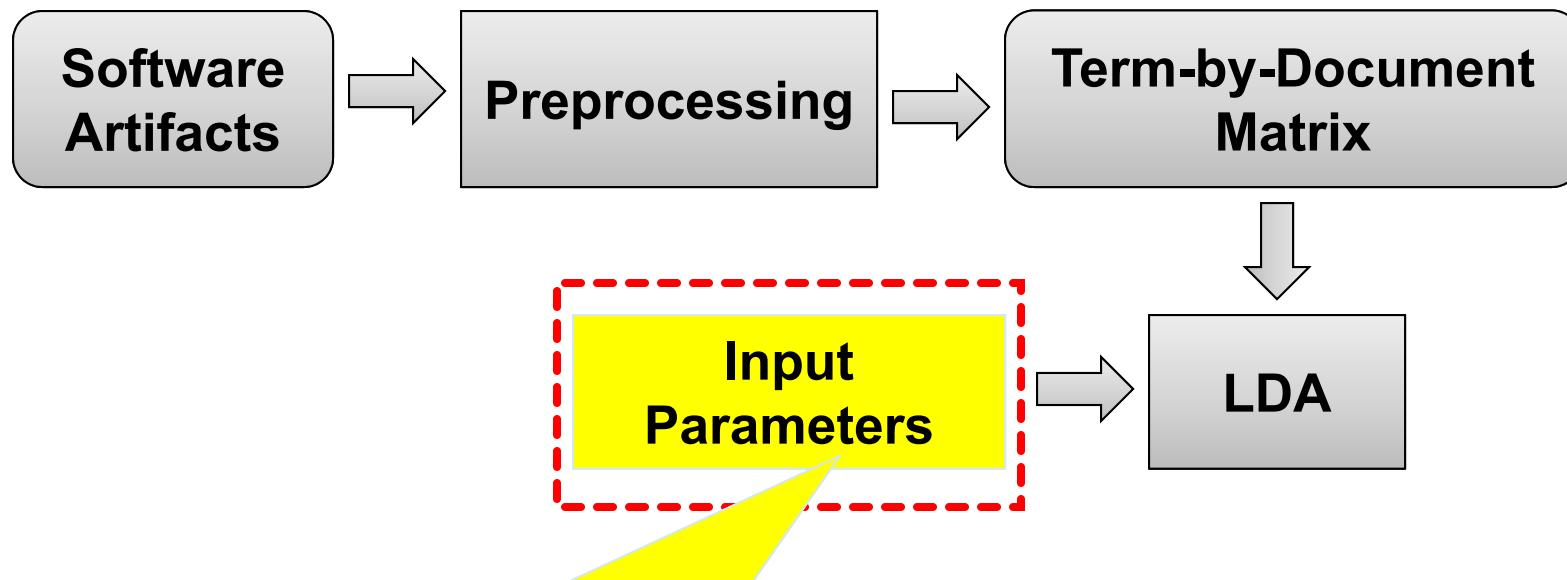






Let's examine the LDA input parameters in more details





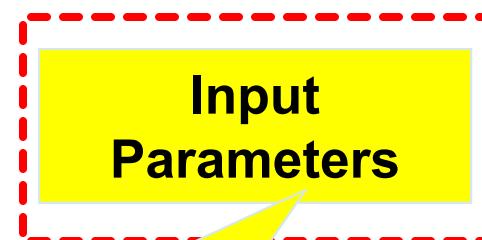
“Configuration”

of iterations

of topics

α

β



“Configuration”

of iterations

of topics

α

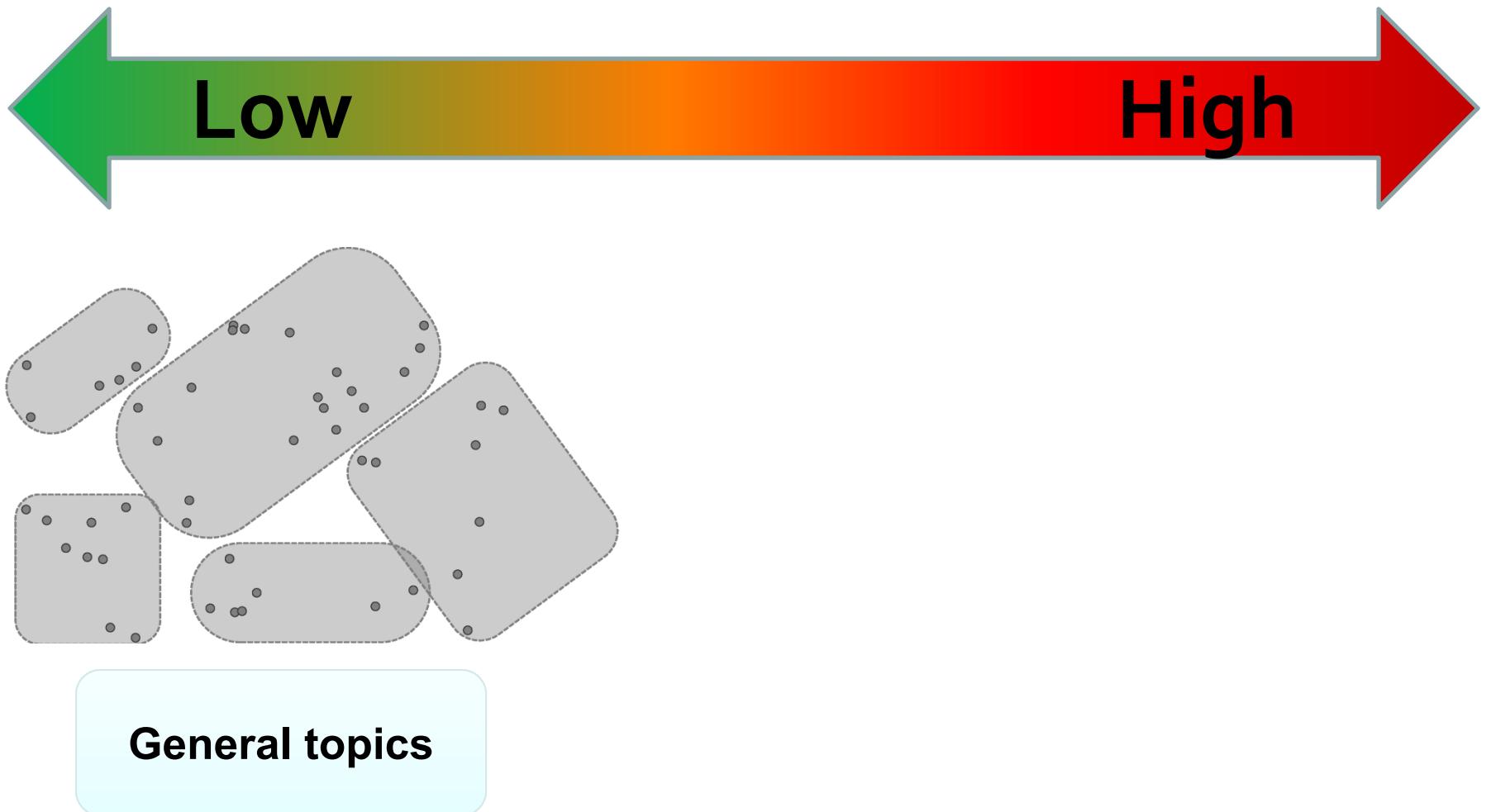
β

Number of Gibbs
samplings of LDA
model

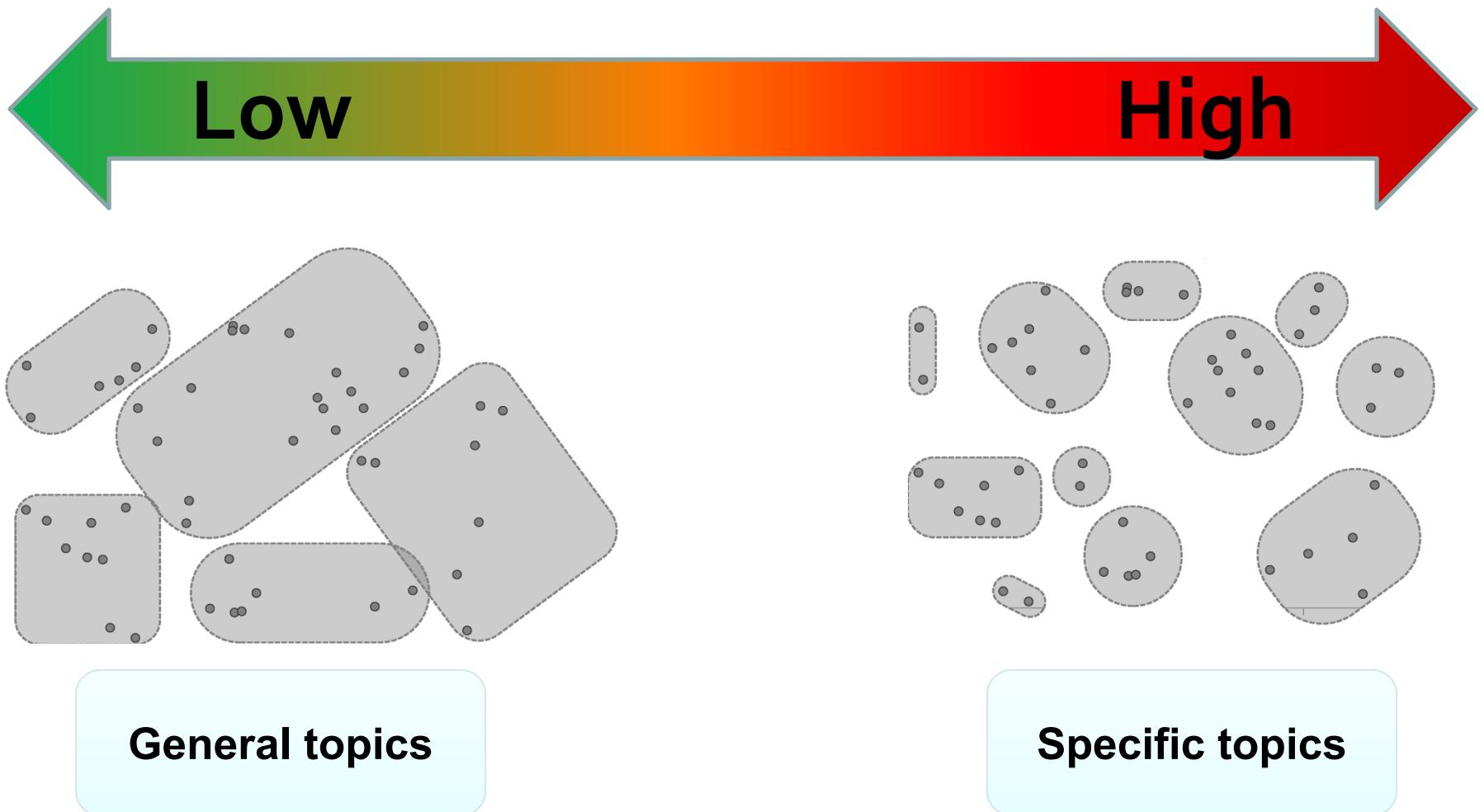
Number of topics...



Number of topics...



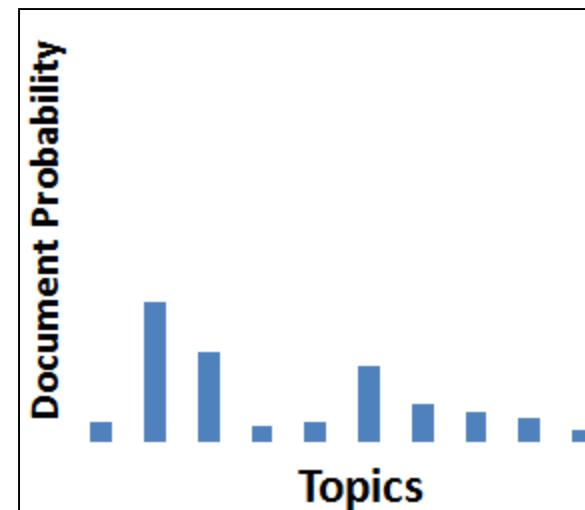
Number of topics...



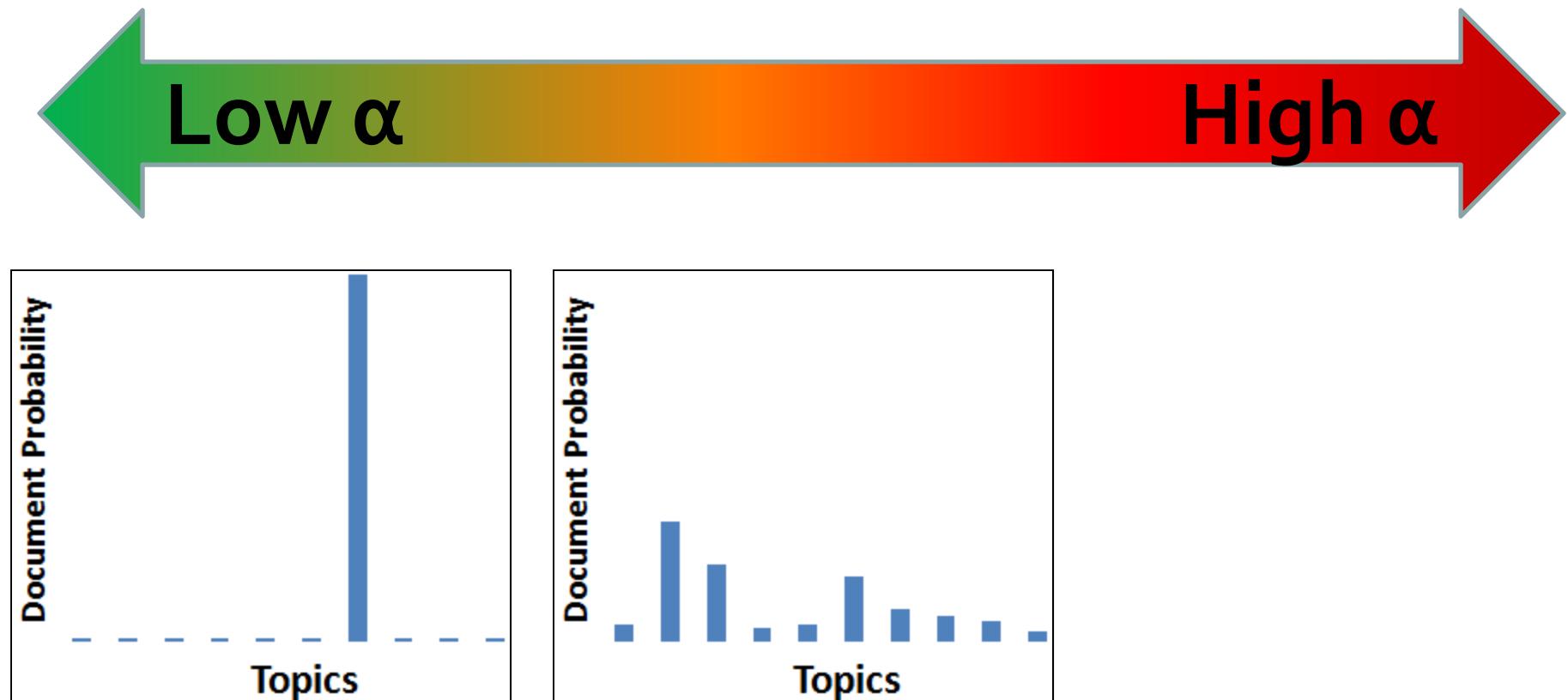
α , influences the “smoothness” of documents to topics distribution



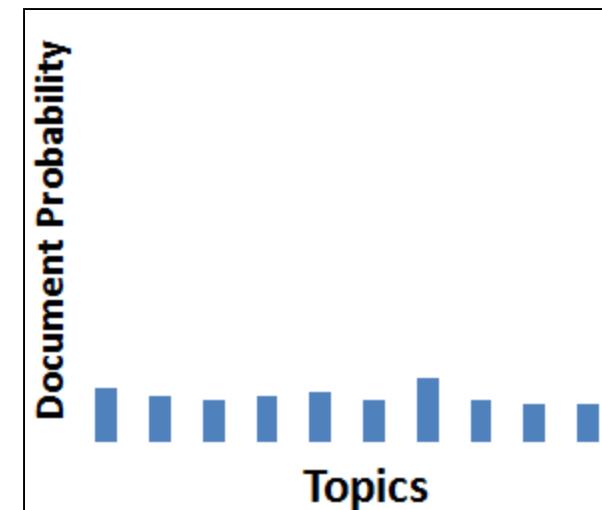
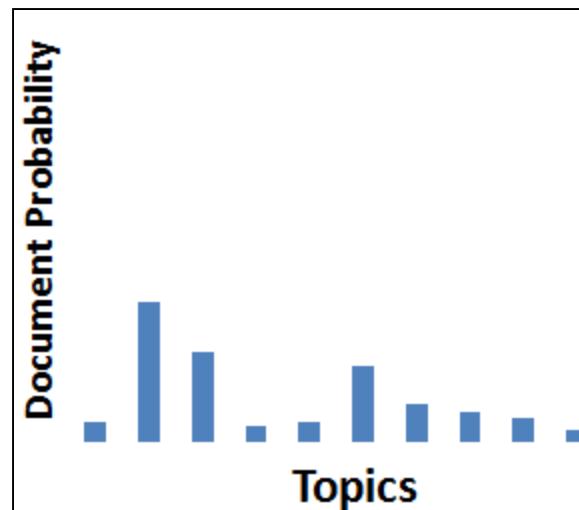
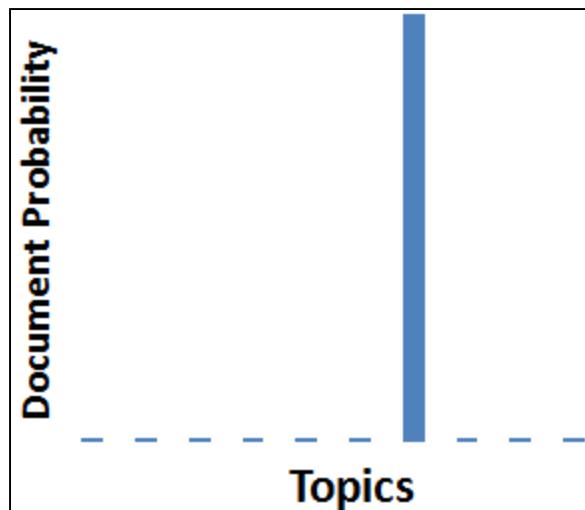
α , influences the “smoothness” of documents to topics distribution



α , influences the “smoothness” of documents to topics distribution



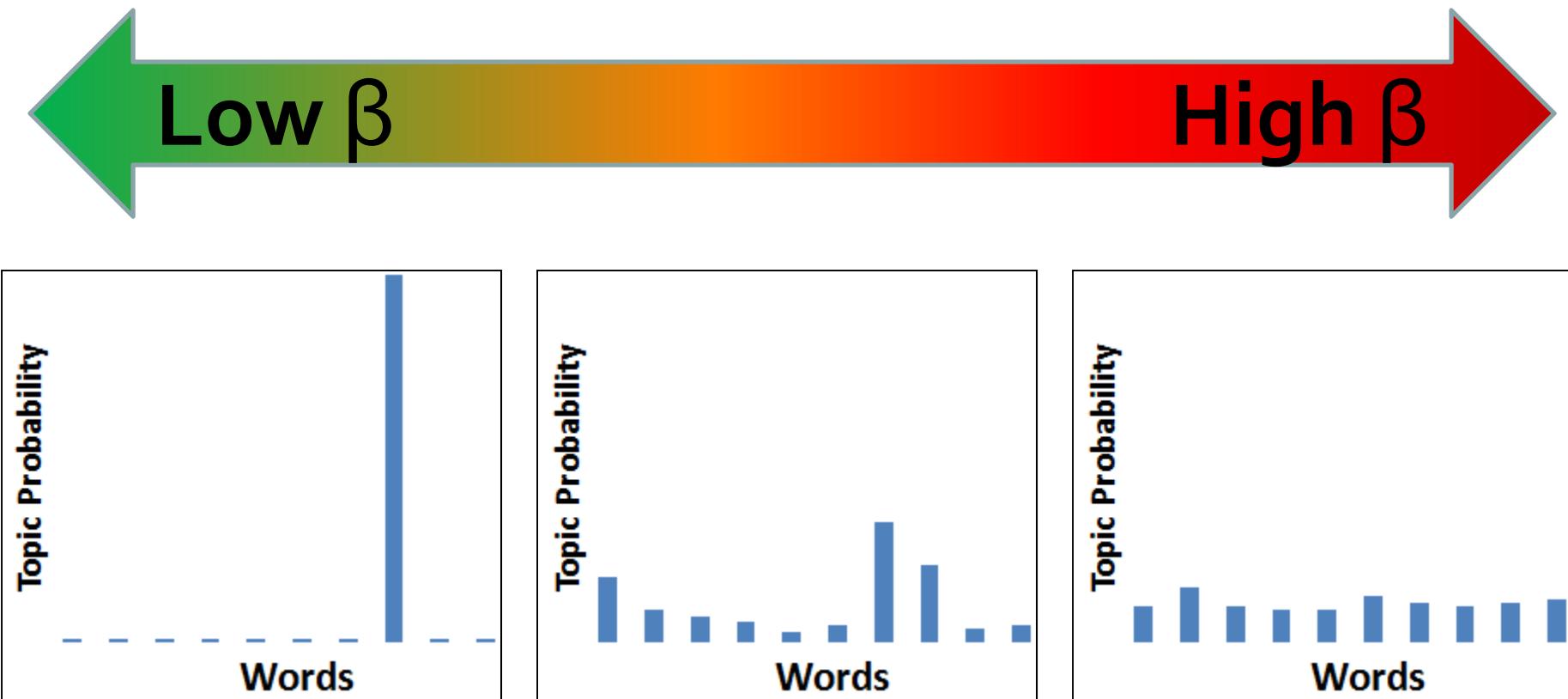
α , influences the “smoothness” of documents to topics distribution



β , influences the “smoothness” of topics to words distribution



β , influences the “smoothness” of topics to words distribution

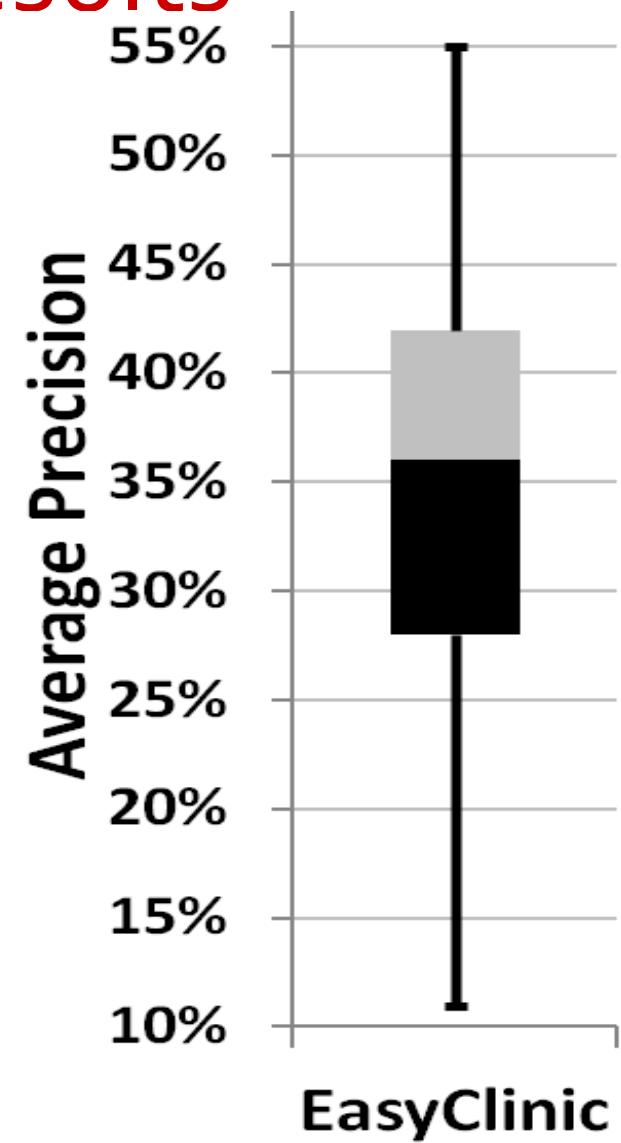


LDA parameters significantly influence the results

- Traceability Link Recovery
- 1,000 different configurations of LDA parameters
 - Evaluate the Average Precision on EasyClinic

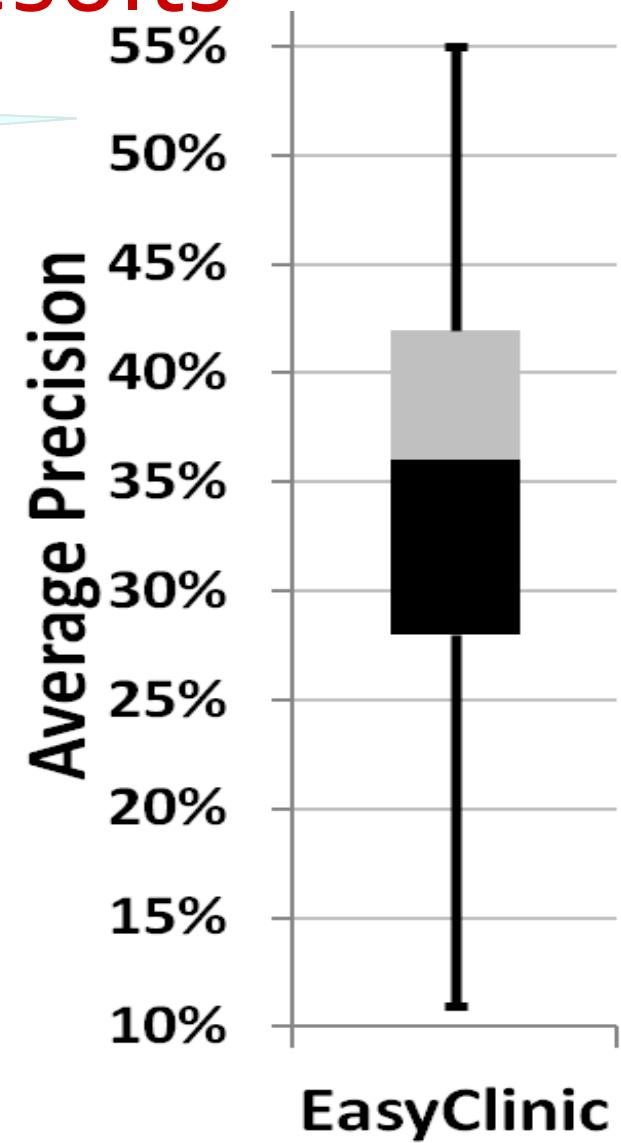
LDA parameters significantly influence the results

- Traceability Link Recovery
- 1,000 different configurations of LDA parameters
 - Evaluate the Average Precision on EasyClinic



LDA parameters significantly influence the results

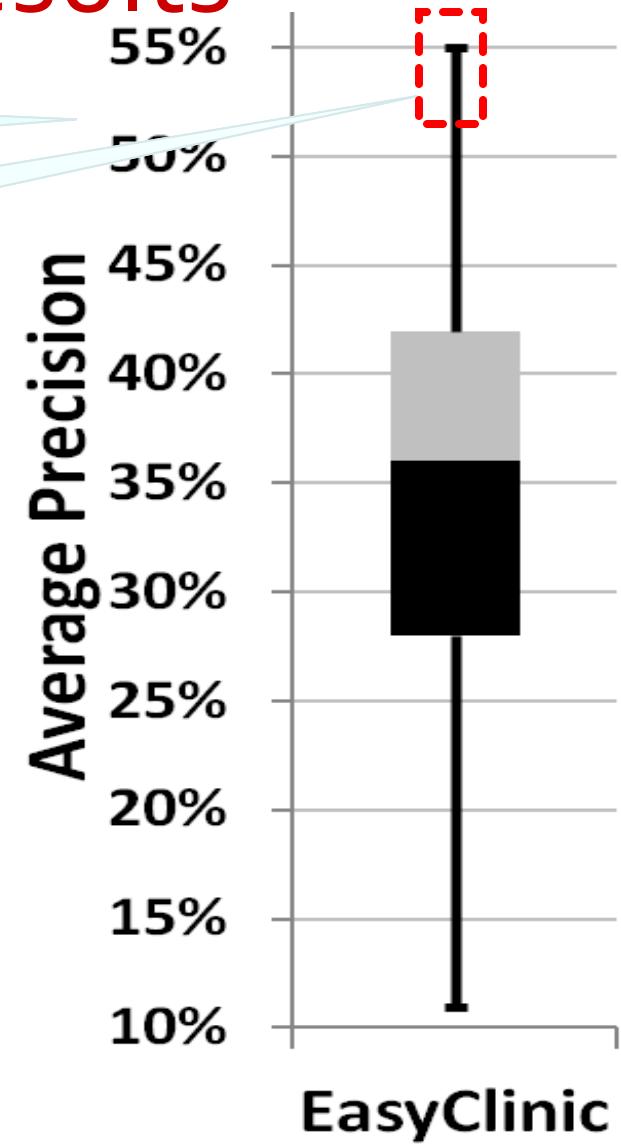
High variability in results



LDA parameters significantly influence the results

High variability in results

Few configurations produce good results



What kind of LDA configurations were used for software?



What kind of LDA configurations were used for software?



“ad-hoc”
configurations

Parameters “imported”
from natural language
community

Assumption:
source code
has the same characteristics as
natural language

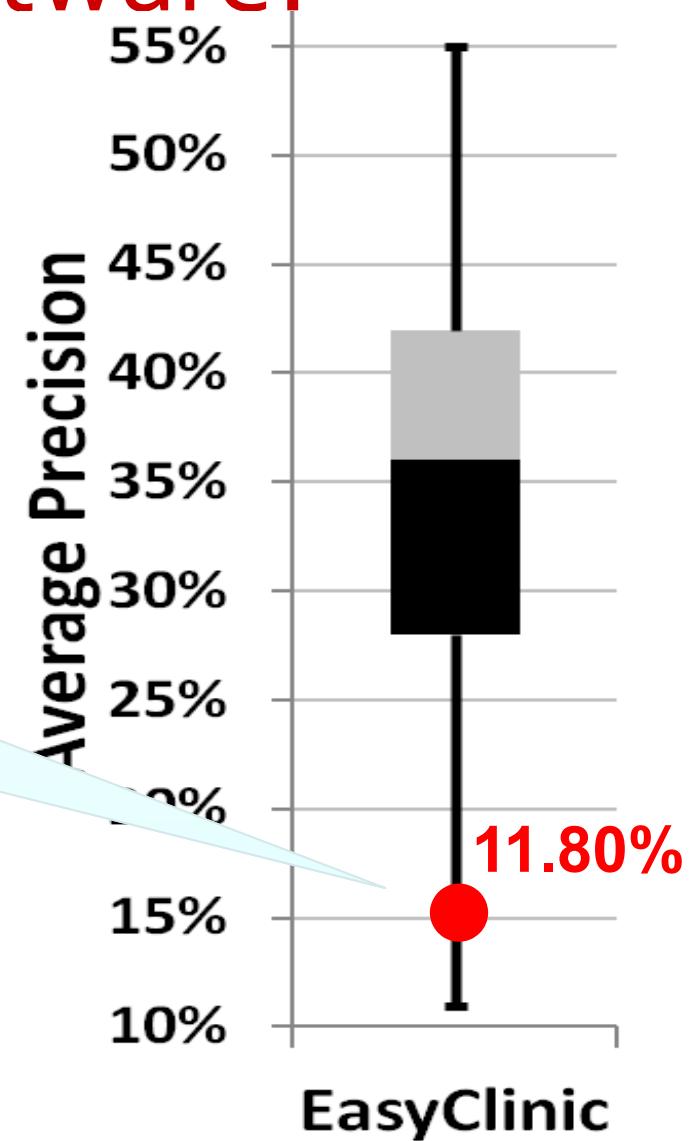
What kind of LDA configurations were used for software?



“ad-hoc” configurations

Parameters “imported” from natural language community

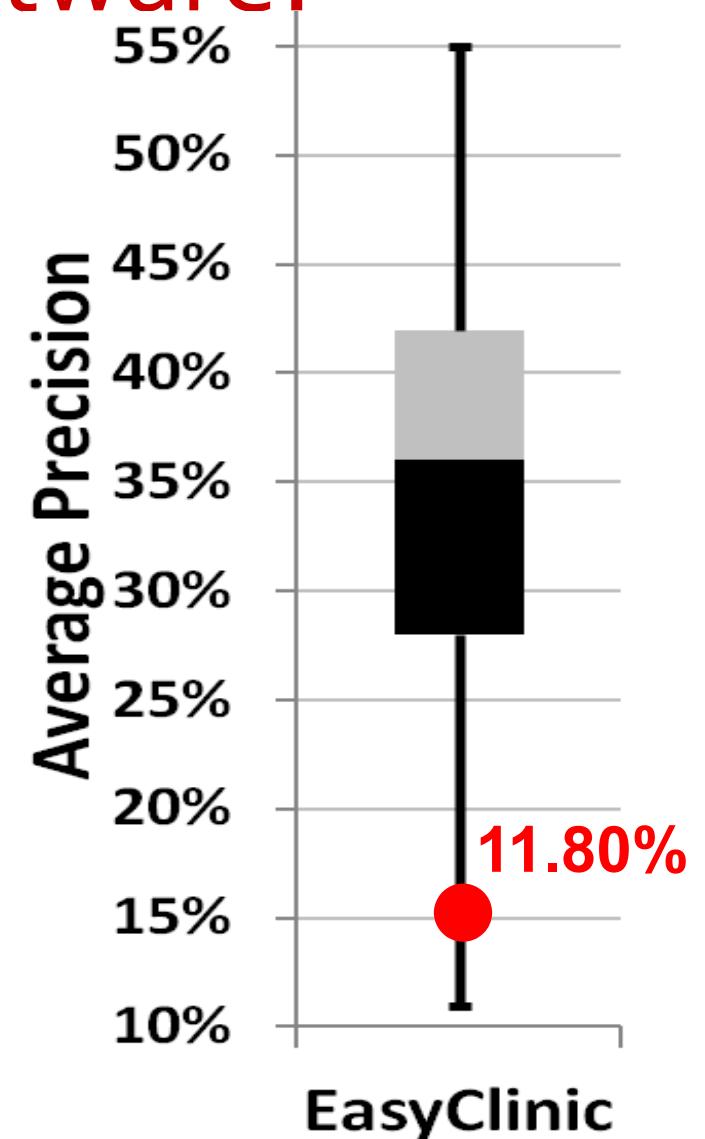
Assumption:
source code
has the same characteristics as
natural language



What kind of LDA configurations were used for software?

[Hindle et al. @ ICSE'12]:
source code
is more *regular* and *predictable* than
natural language

~~Assumption:~~
~~source code~~
has the same characteristics as
natural language

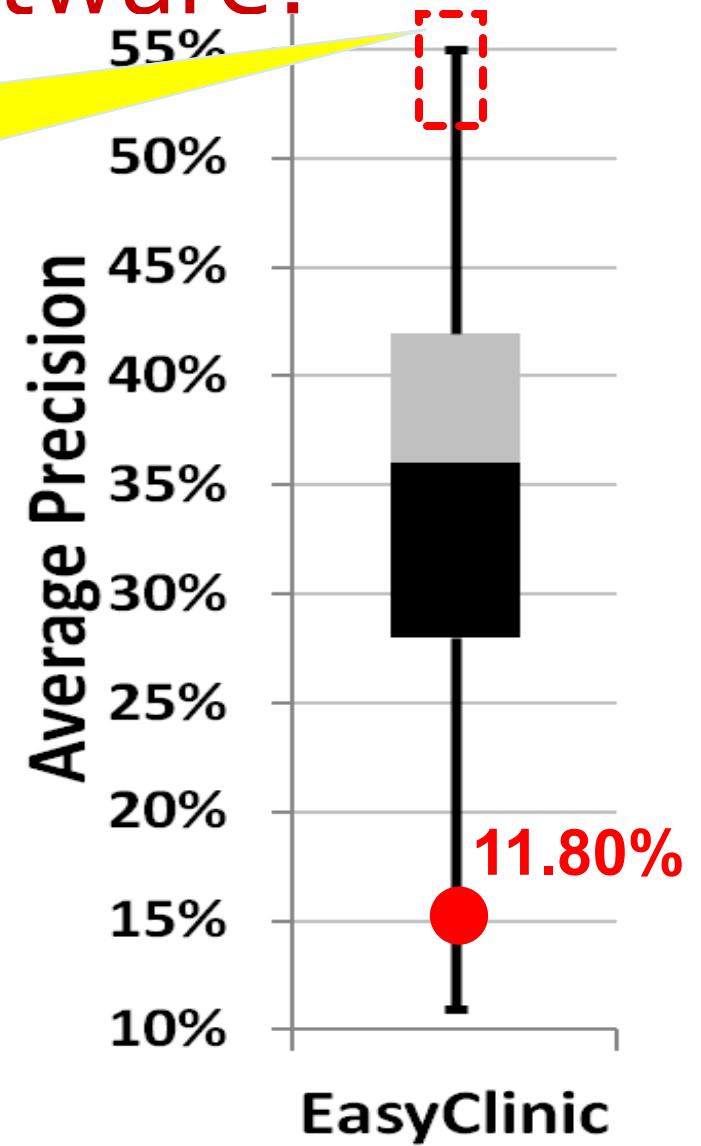


What kind of LDA configurations were used for software?

We need new techniques to find these configurations

[Hindle et al. @ ICSE'12]:
source code
is more *regular* and *predictable* than
natural language

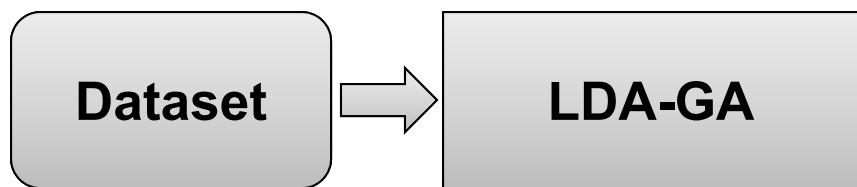
~~Assumption:~~
~~source code~~
has the same characteristics as
natural language



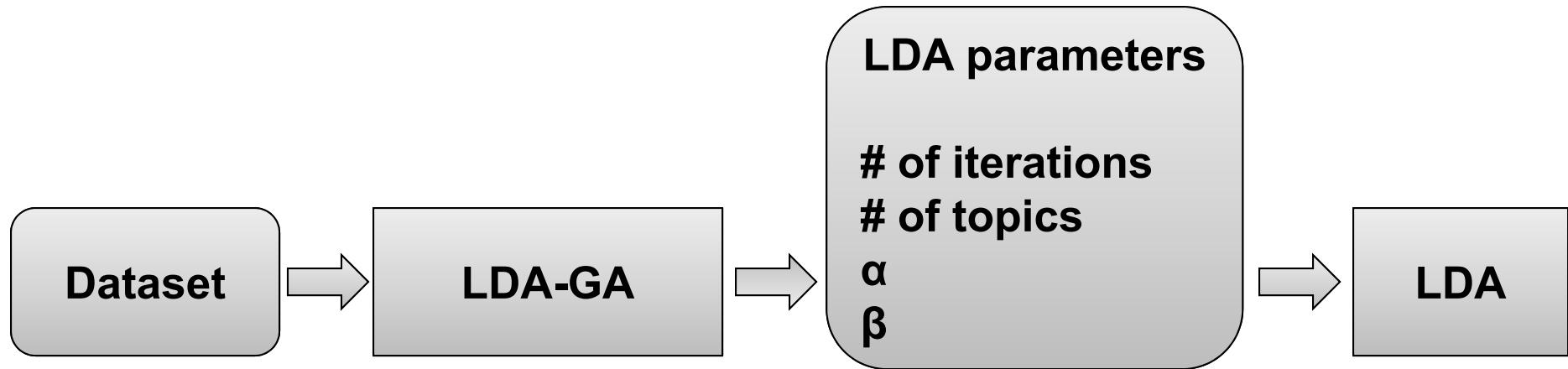
Our contribution...**LDA-GA**

LDA-GA: automatically calibrate the input parameters of LDA using a genetic algorithm

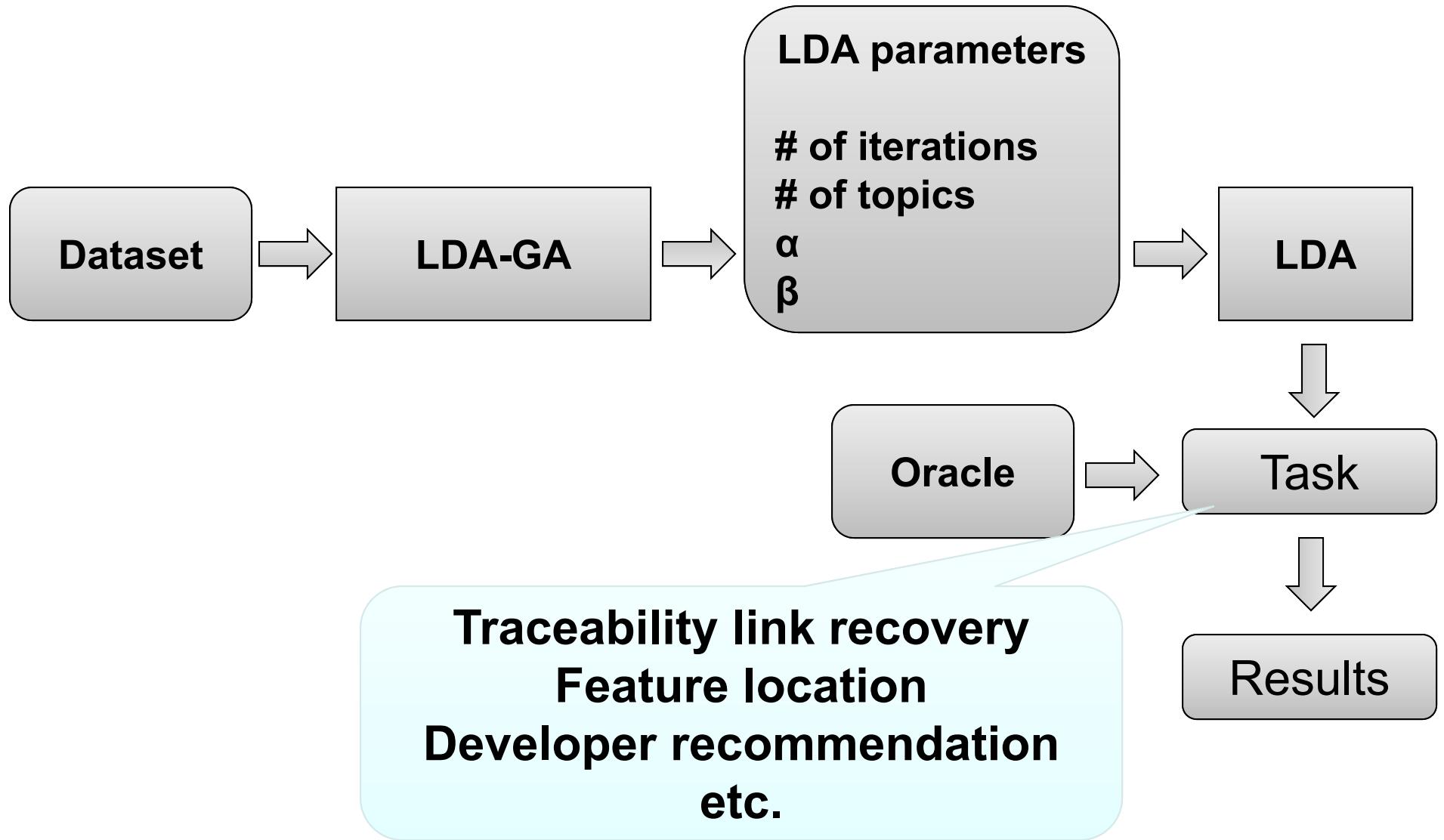
LDA-GA: automatically calibrate the input parameters of LDA using a genetic algorithm



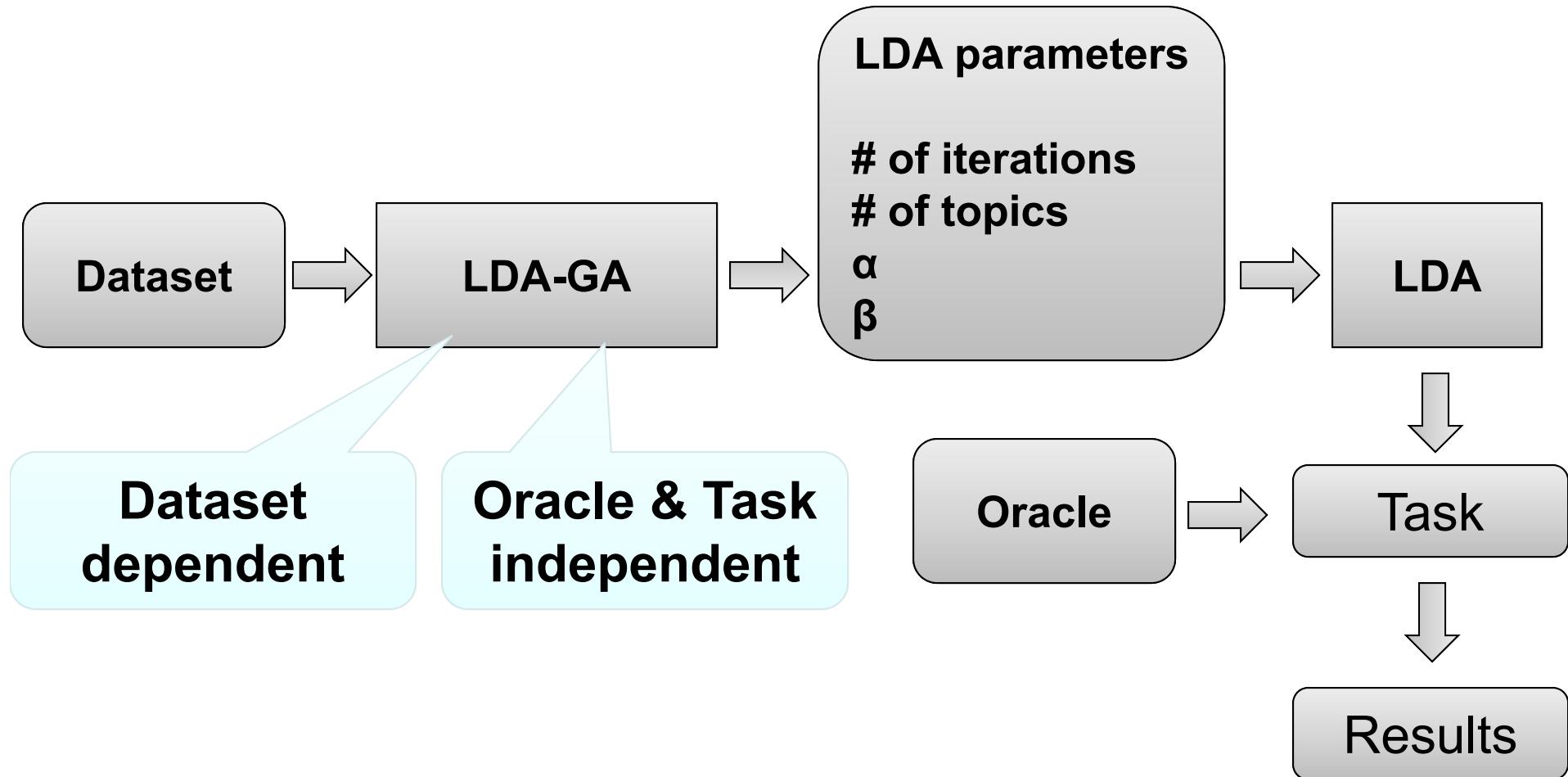
LDA-GA: automatically calibrate the input parameters of LDA using a genetic algorithm



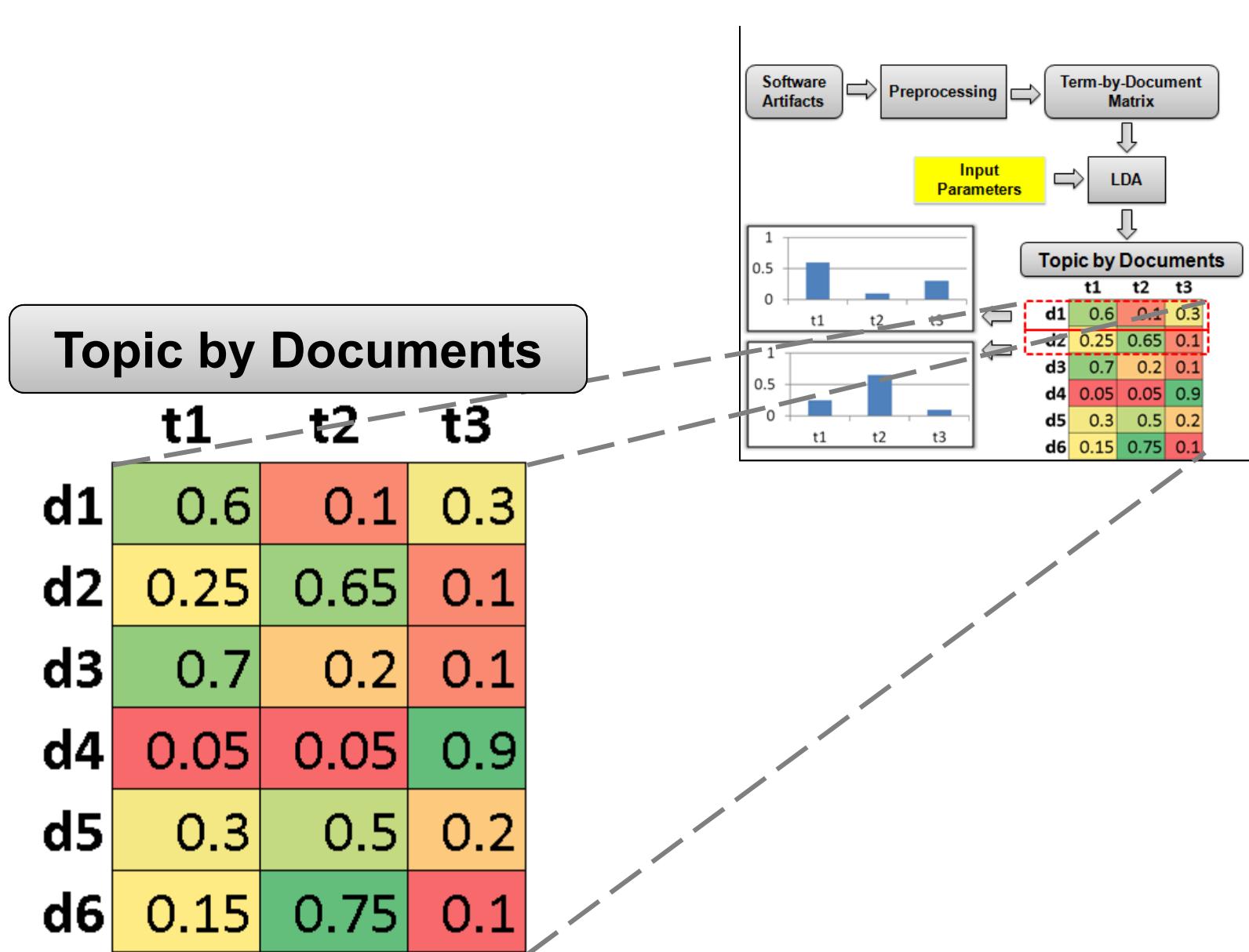
LDA-GA: automatically calibrate the input parameters of LDA using a genetic algorithm



LDA-GA: automatically calibrate the input parameters of LDA using a genetic algorithm



**How to evaluate how “good” an
LDA configuration is?**



Topic by Documents

| | t1 | t2 | t3 |
|----|------|------|-----|
| d1 | 0.6 | 0.1 | 0.3 |
| d2 | 0.25 | 0.65 | 0.1 |
| d3 | 0.7 | 0.2 | 0.1 |
| d4 | 0.05 | 0.05 | 0.9 |
| d5 | 0.3 | 0.5 | 0.2 |
| d6 | 0.15 | 0.75 | 0.1 |

Topic by Documents

| | t1 | t2 | t3 |
|----|------|------|-----|
| d1 | 0.6 | 0.1 | 0.3 |
| d2 | 0.25 | 0.65 | 0.1 |
| d3 | 0.7 | 0.2 | 0.1 |
| d4 | 0.05 | 0.05 | 0.9 |
| d5 | 0.3 | 0.5 | 0.2 |
| d6 | 0.15 | 0.75 | 0.1 |

Dominant
Topics

Topic by Documents

| | t1 | t2 | t3 |
|----|------|------|-----|
| d1 | 0.6 | 0.1 | 0.3 |
| d2 | 0.25 | 0.65 | 0.1 |
| d3 | 0.7 | 0.2 | 0.1 |
| d4 | 0.05 | 0.05 | 0.9 |
| d5 | 0.3 | 0.5 | 0.2 |
| d6 | 0.15 | 0.75 | 0.1 |

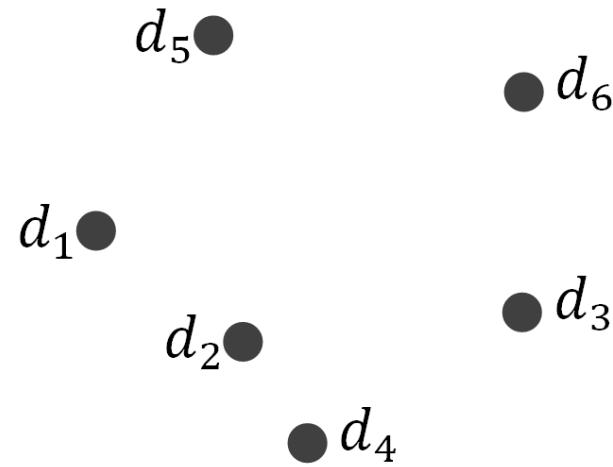
Dominant
Topics

Topic by Documents

| | t1 | t2 | t3 |
|----|------|------|-----|
| d1 | 0.6 | 0.1 | 0.3 |
| d2 | 0.25 | 0.65 | 0.1 |
| d3 | 0.7 | 0.2 | 0.1 |
| d4 | 0.05 | 0.05 | 0.9 |
| d5 | 0.3 | 0.5 | 0.2 |
| d6 | 0.15 | 0.75 | 0.1 |

Dominant
Topics

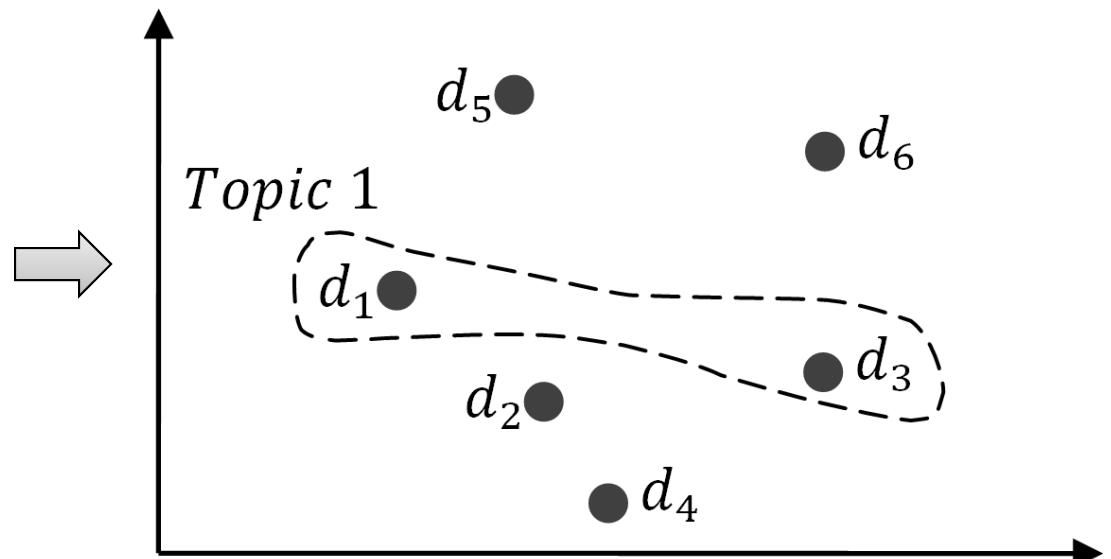
LDA Model



Topic by Documents

| | t1 | t2 | t3 |
|----|------|------|-----|
| d1 | 0.6 | 0.1 | 0.3 |
| d2 | 0.25 | 0.65 | 0.1 |
| d3 | 0.7 | 0.2 | 0.1 |
| d4 | 0.05 | 0.05 | 0.9 |
| d5 | 0.3 | 0.5 | 0.2 |
| d6 | 0.15 | 0.75 | 0.1 |

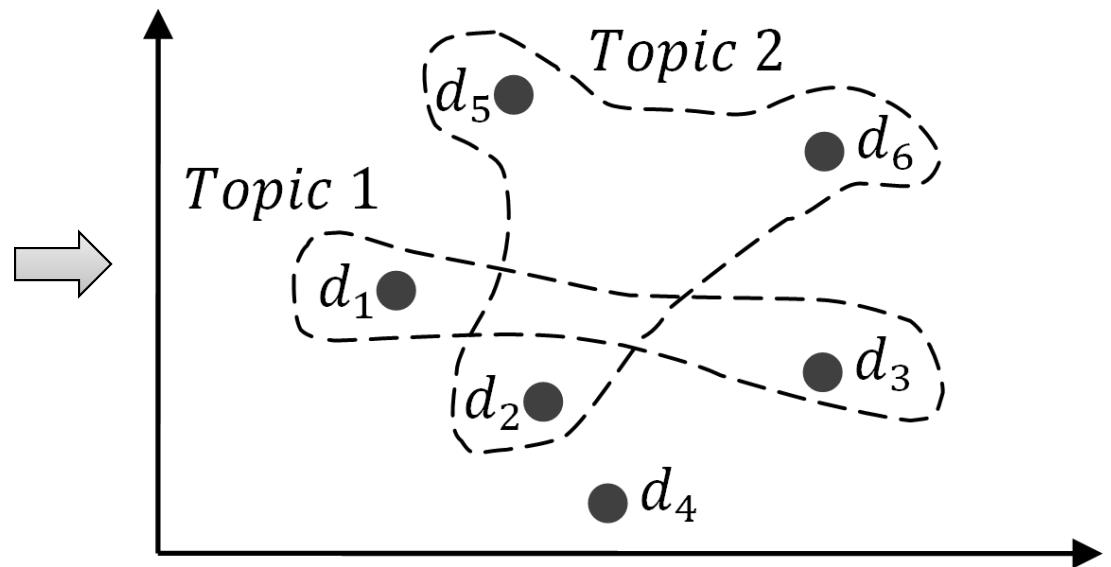
LDA Model



Topic by Documents

| | t1 | t2 | t3 |
|----|------|------|-----|
| d1 | 0.6 | 0.1 | 0.3 |
| d2 | 0.25 | 0.65 | 0.1 |
| d3 | 0.7 | 0.2 | 0.1 |
| d4 | 0.05 | 0.05 | 0.9 |
| d5 | 0.3 | 0.5 | 0.2 |
| d6 | 0.15 | 0.75 | 0.1 |

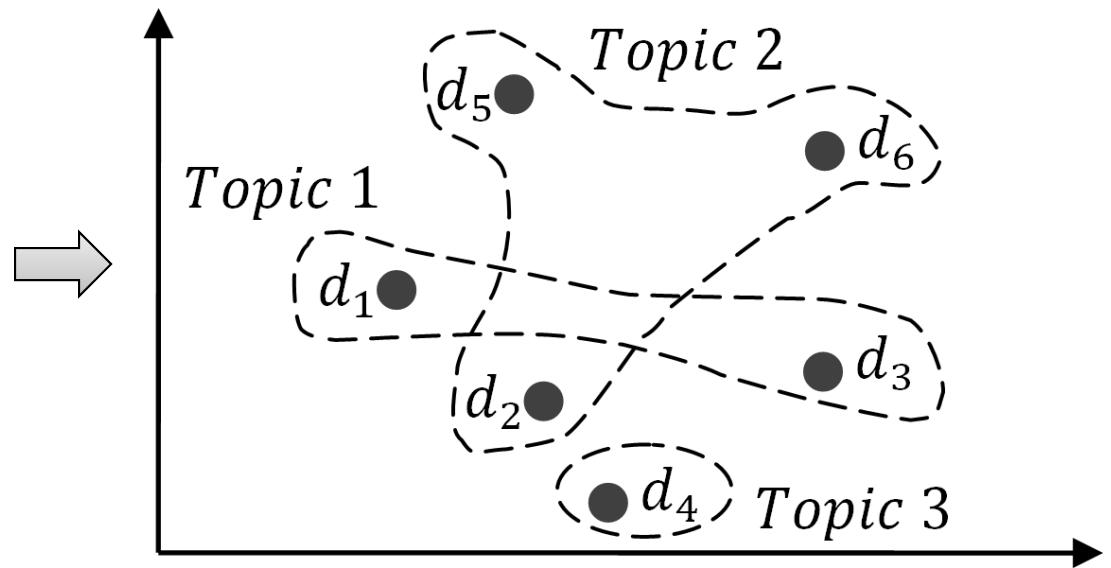
LDA Model



Topic by Documents

| | t1 | t2 | t3 |
|----|------|------|-----|
| d1 | 0.6 | 0.1 | 0.3 |
| d2 | 0.25 | 0.65 | 0.1 |
| d3 | 0.7 | 0.2 | 0.1 |
| d4 | 0.05 | 0.05 | 0.9 |
| d5 | 0.3 | 0.5 | 0.2 |
| d6 | 0.15 | 0.75 | 0.1 |

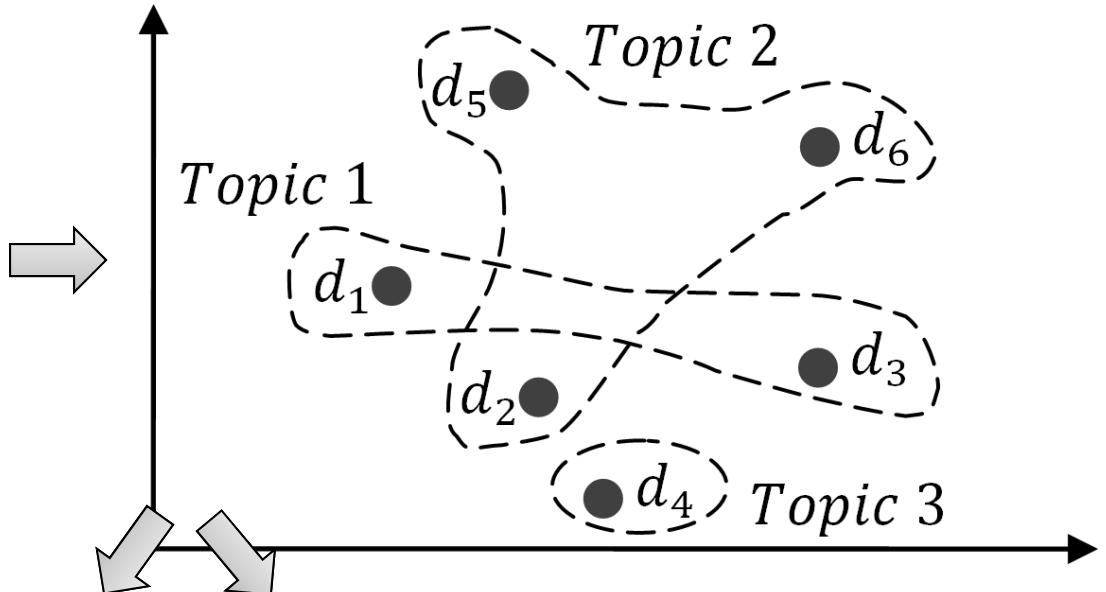
LDA Model



Topic by Documents

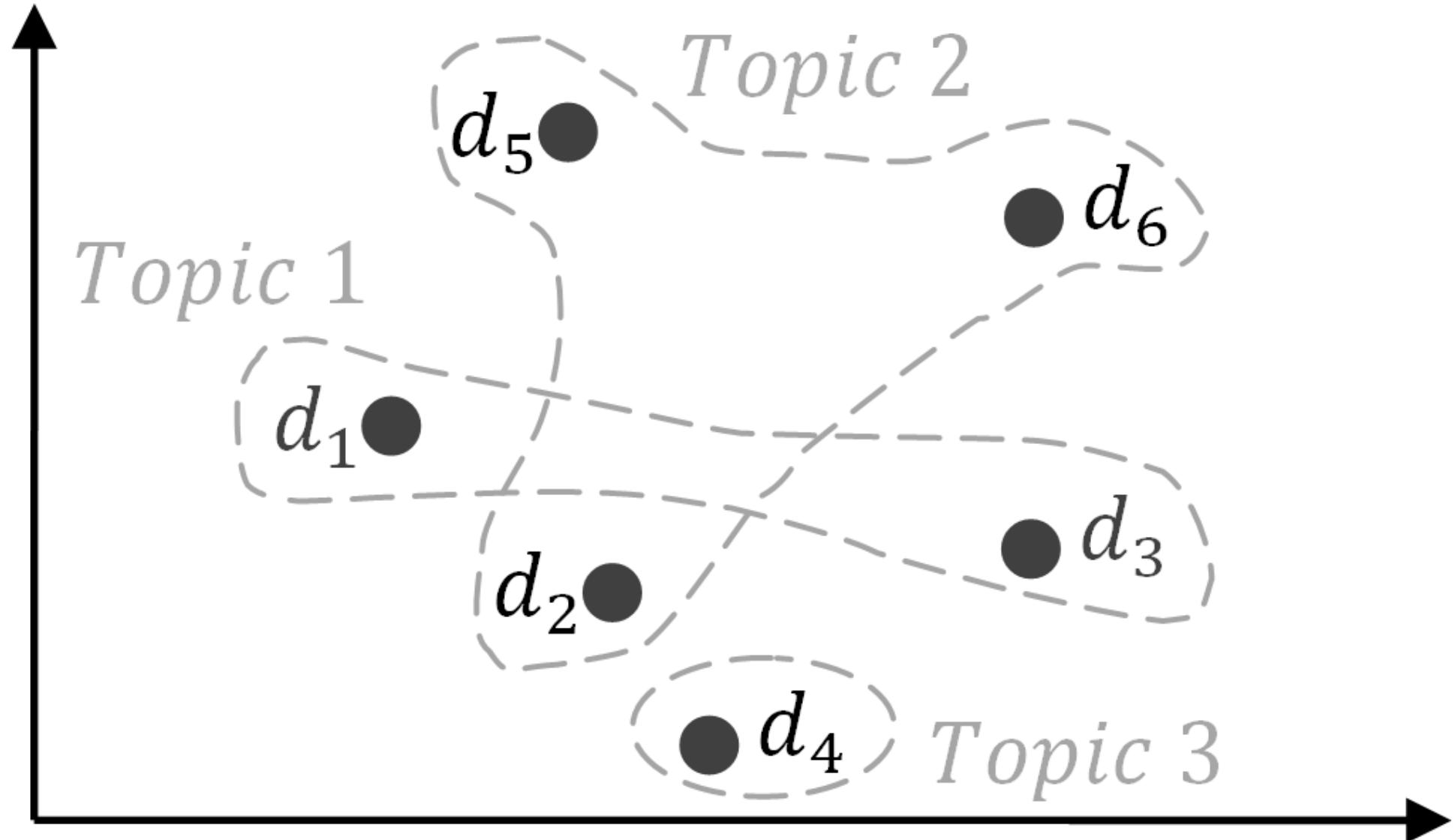
| | t1 | t2 | t3 |
|----|------|------|-----|
| d1 | 0.6 | 0.1 | 0.3 |
| d2 | 0.25 | 0.65 | 0.1 |
| d3 | 0.7 | 0.2 | 0.1 |
| d4 | 0.05 | 0.05 | 0.9 |
| d5 | 0.3 | 0.5 | 0.2 |
| d6 | 0.15 | 0.75 | 0.1 |

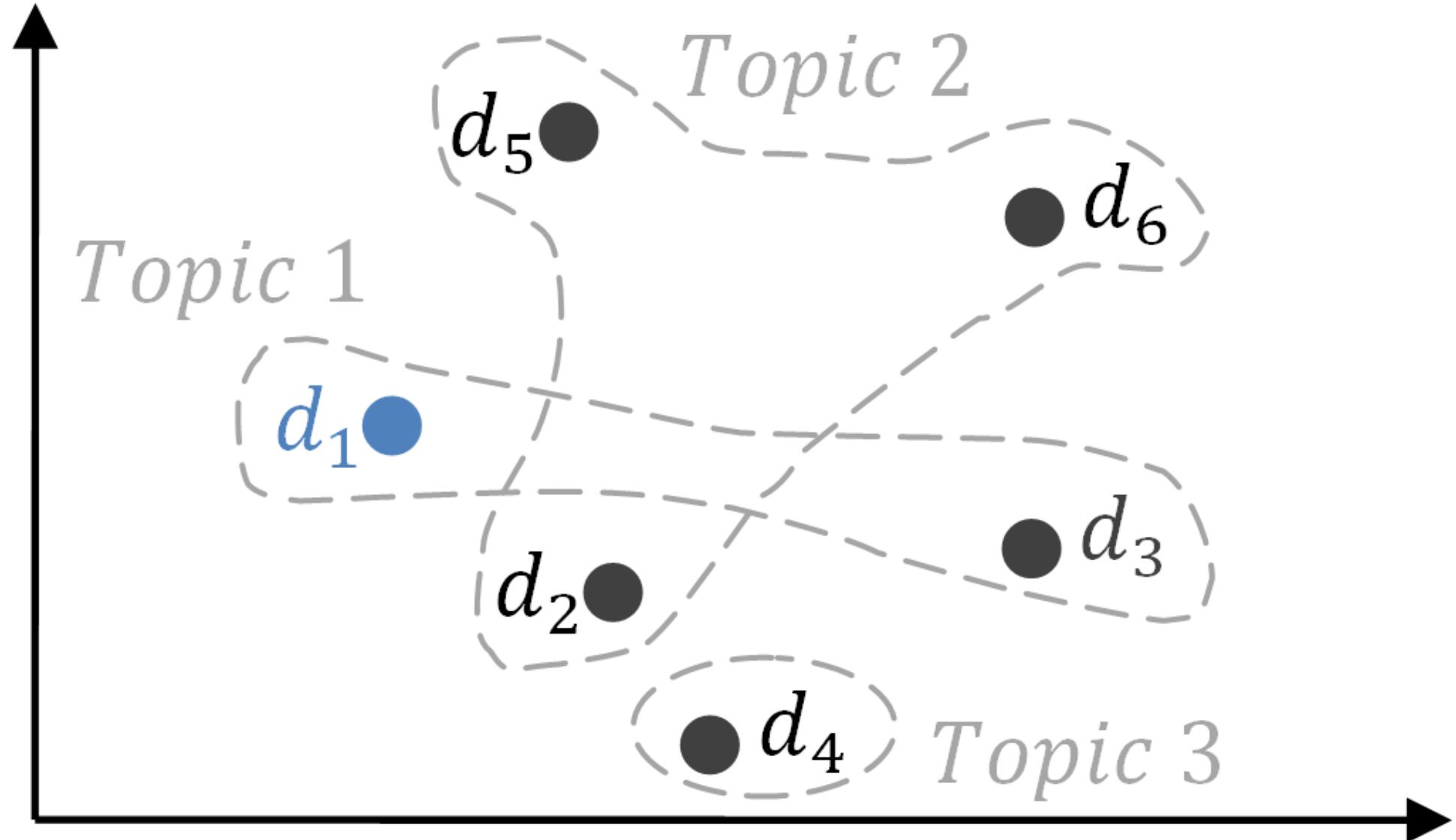
LDA Model

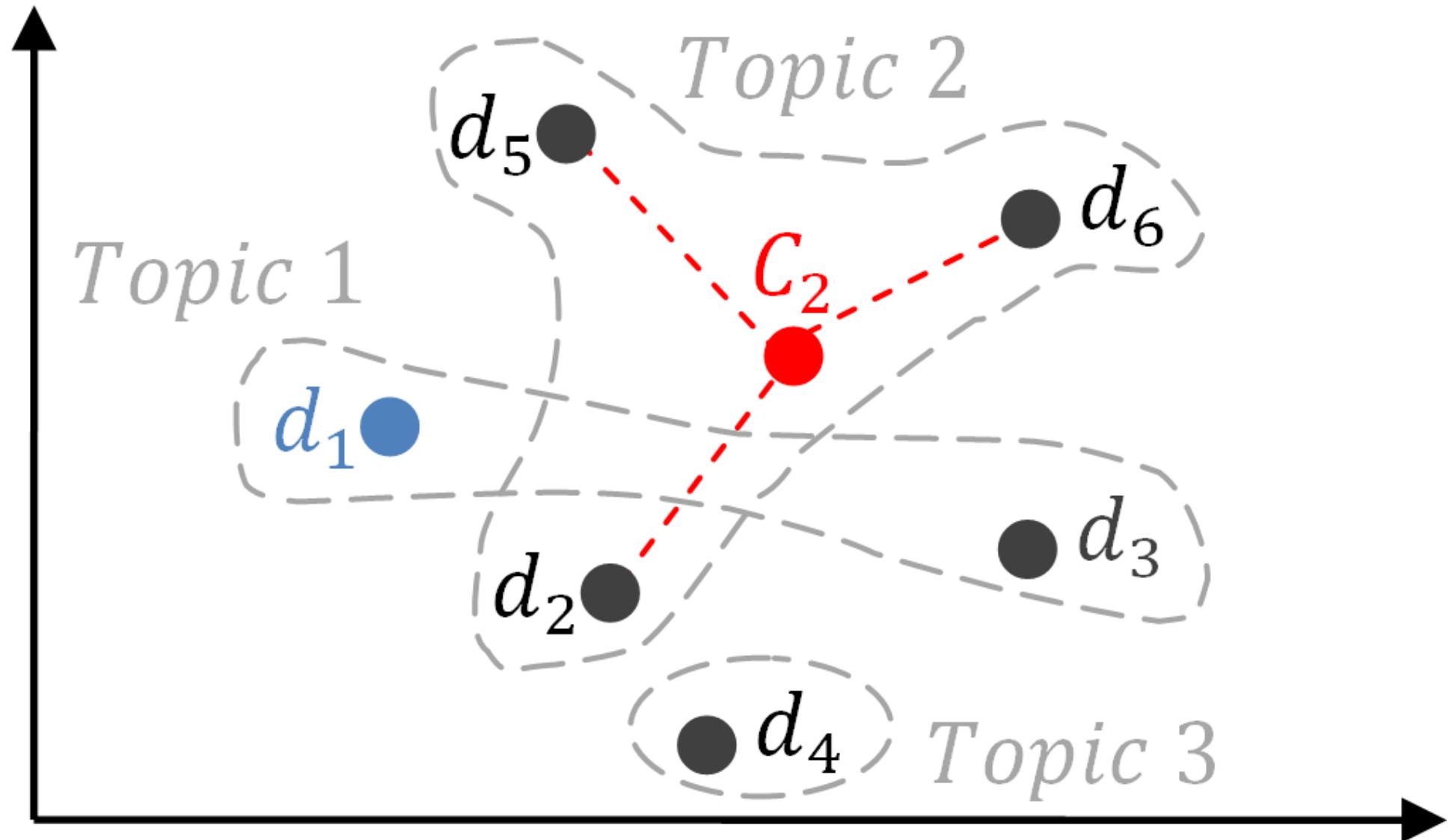


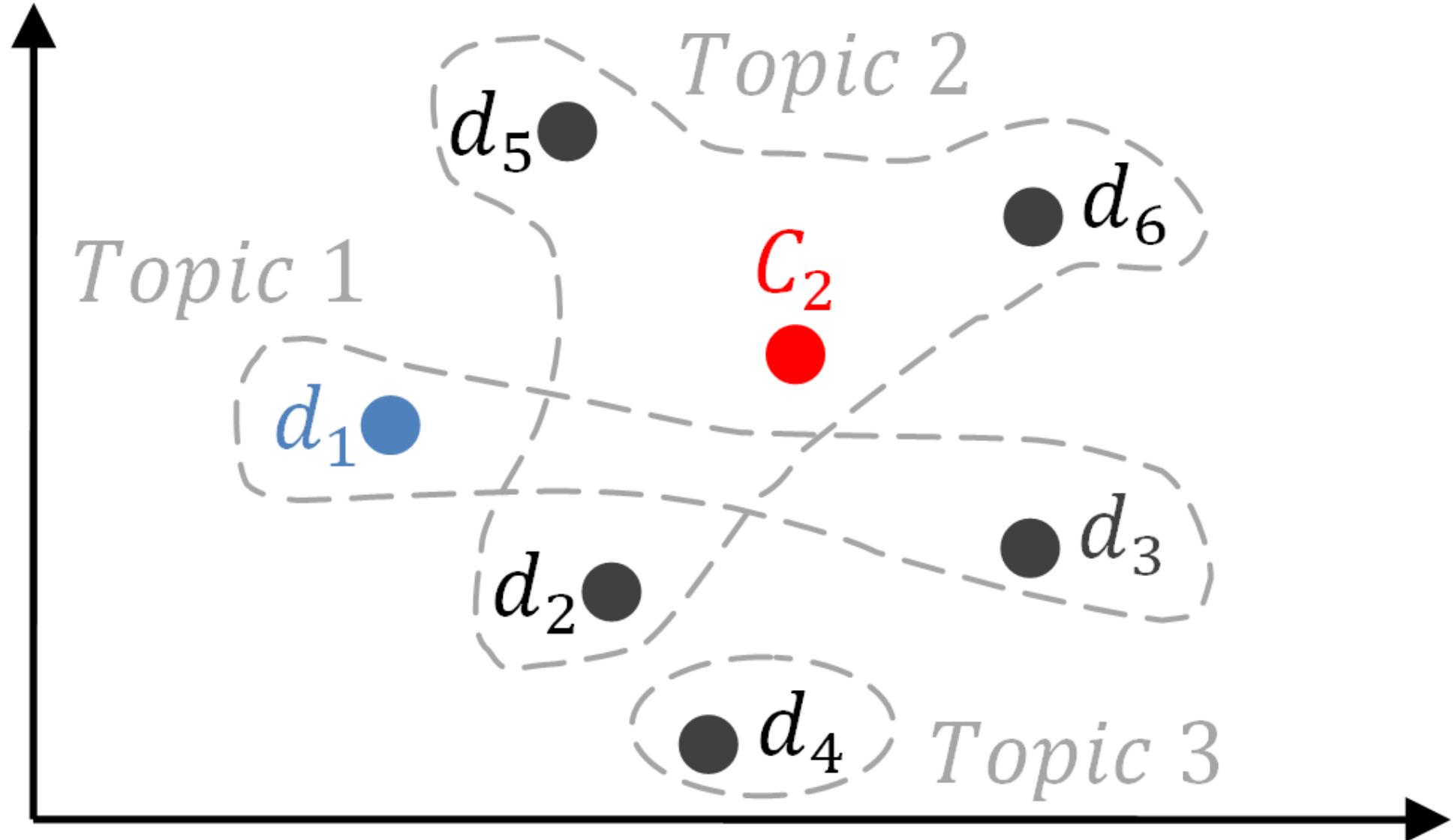
Cohesion (similarity): how related the documents in the same clusters are

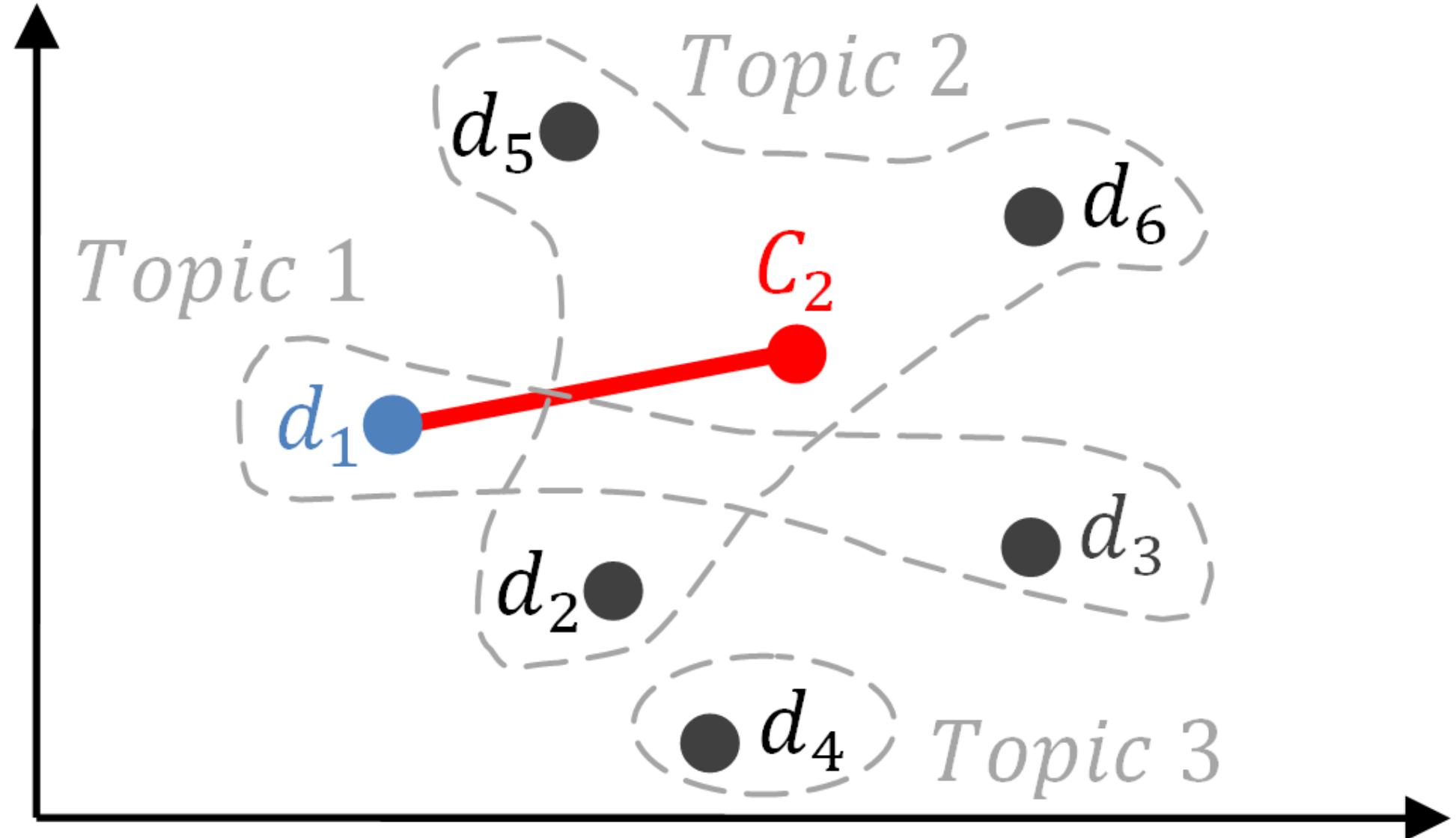
Separation (dissimilarity): how distinct a cluster is from other clusters

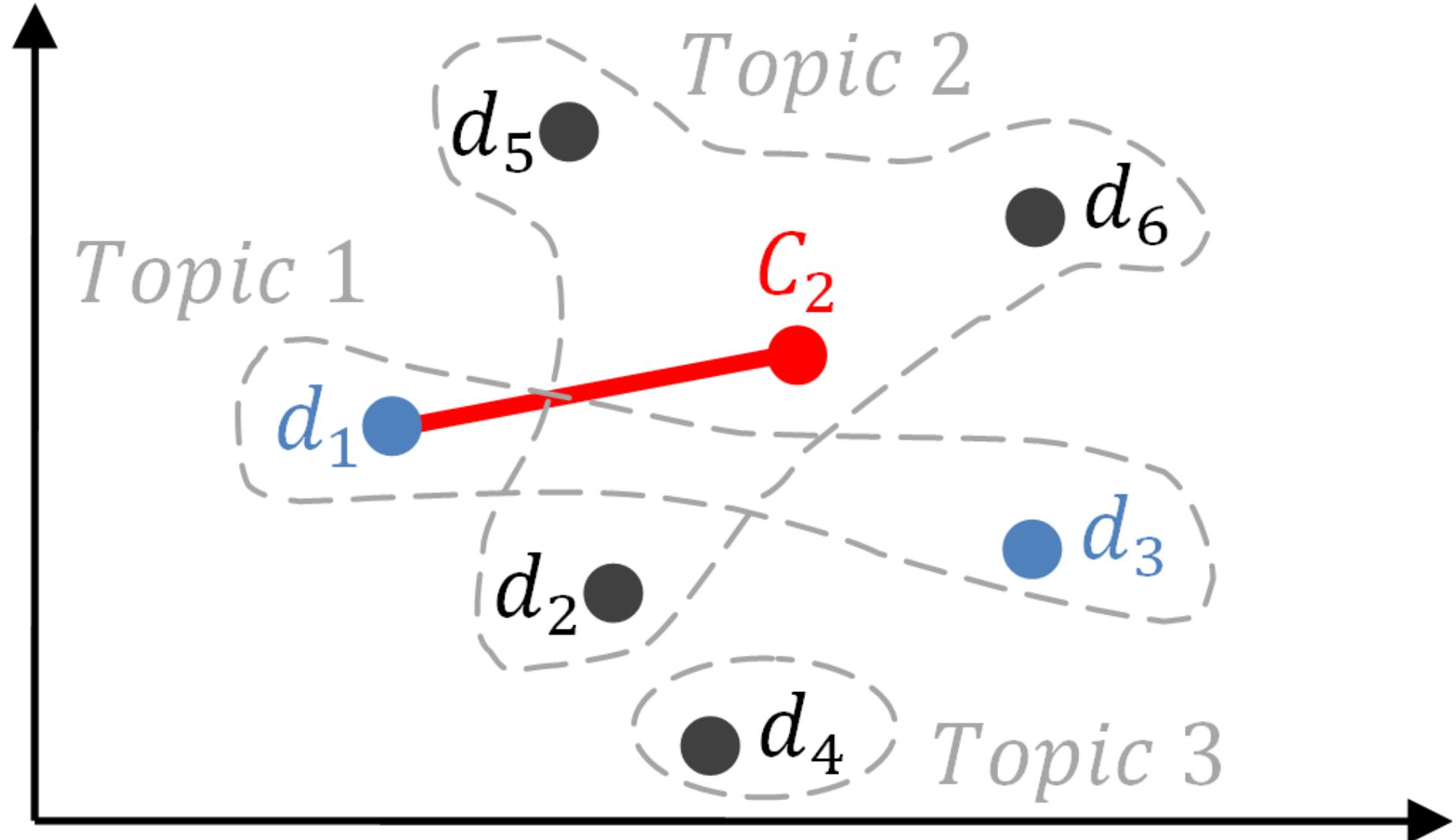


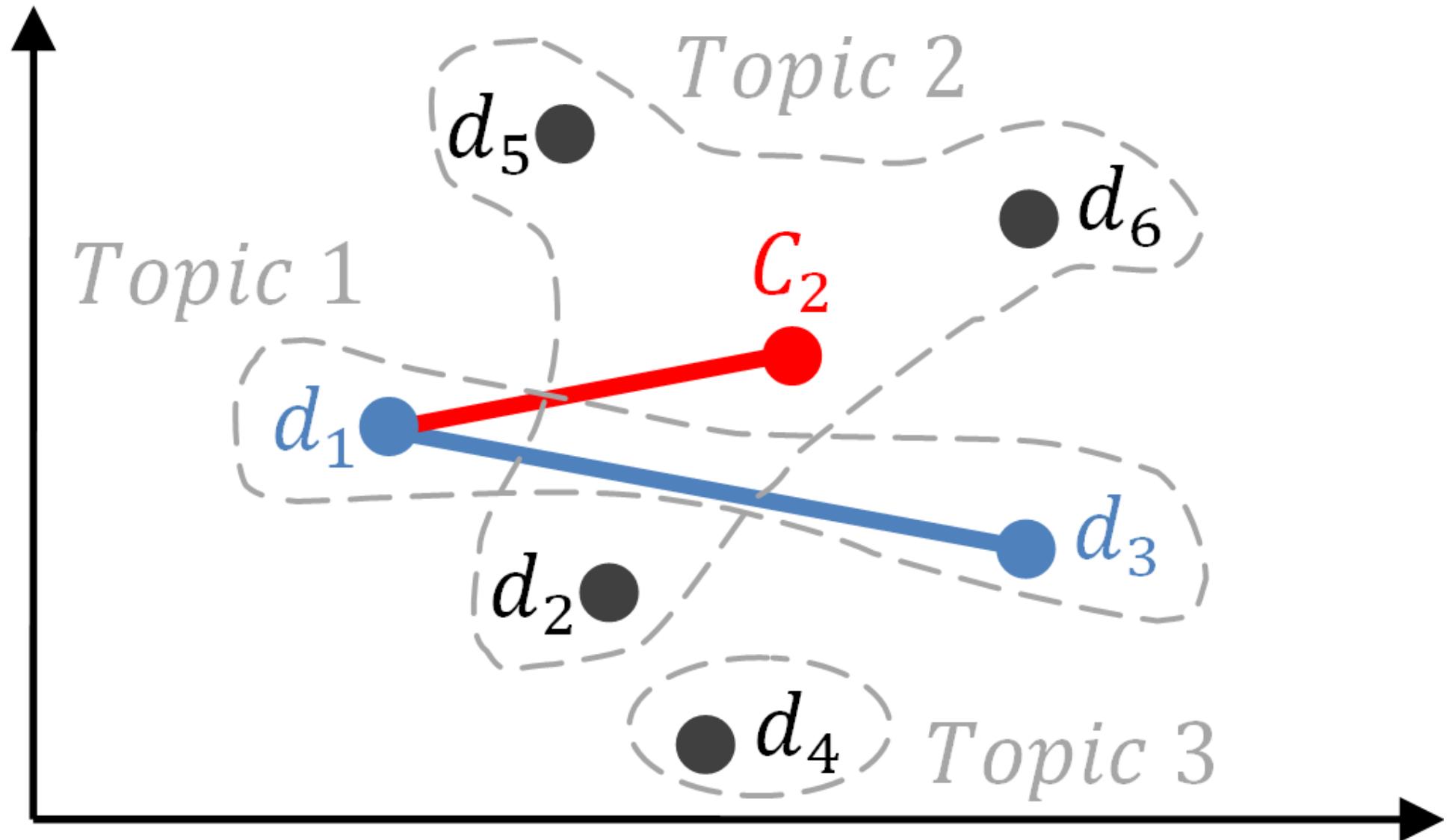


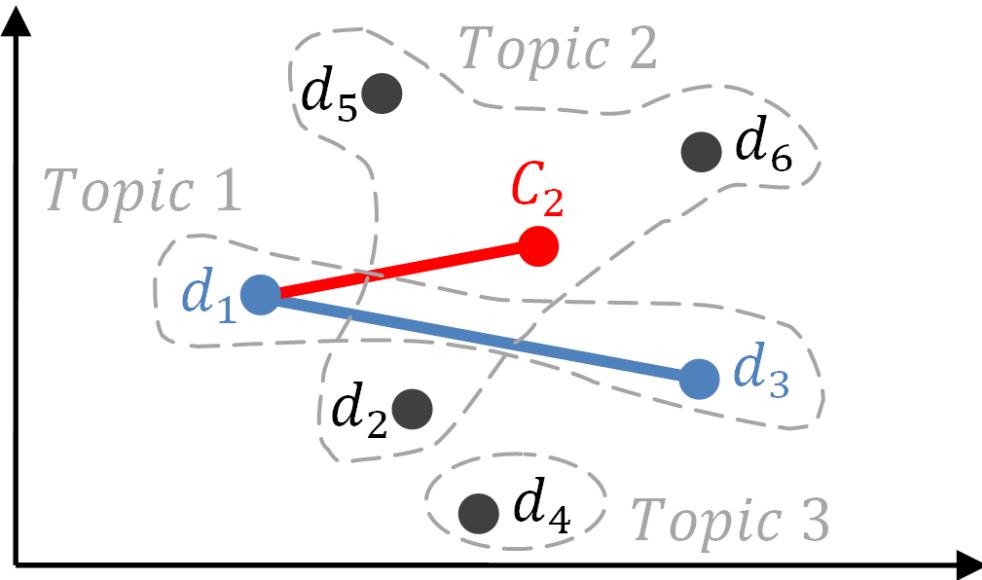




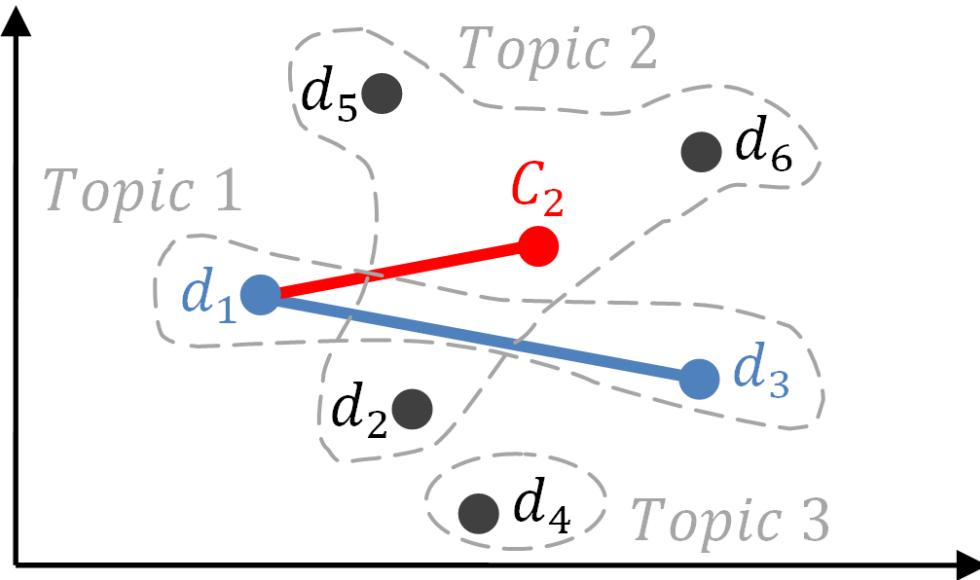








$$silhouette(doc) = \frac{red - blue}{\max(red, blue)}$$



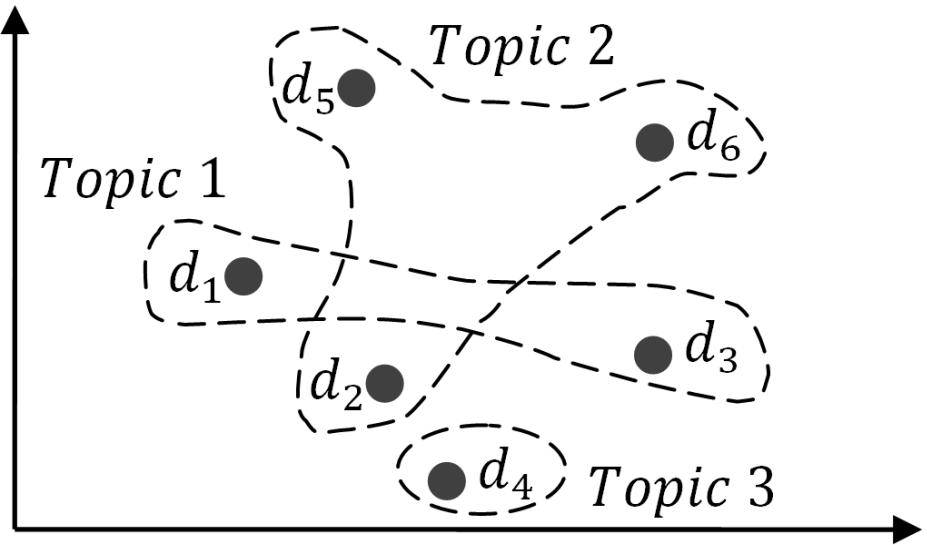
$$\text{silhouette}(doc) = \frac{\text{red} - \text{blue}}{\max(\text{red}, \text{blue})}$$

$$\text{silhouette}(LDA \text{ model}) = \frac{\text{silhouette}(doc_1) + \dots + \text{silhouette}(doc_n)}{\text{number of documents}}$$

**Higher silhouette
coefficients are better**

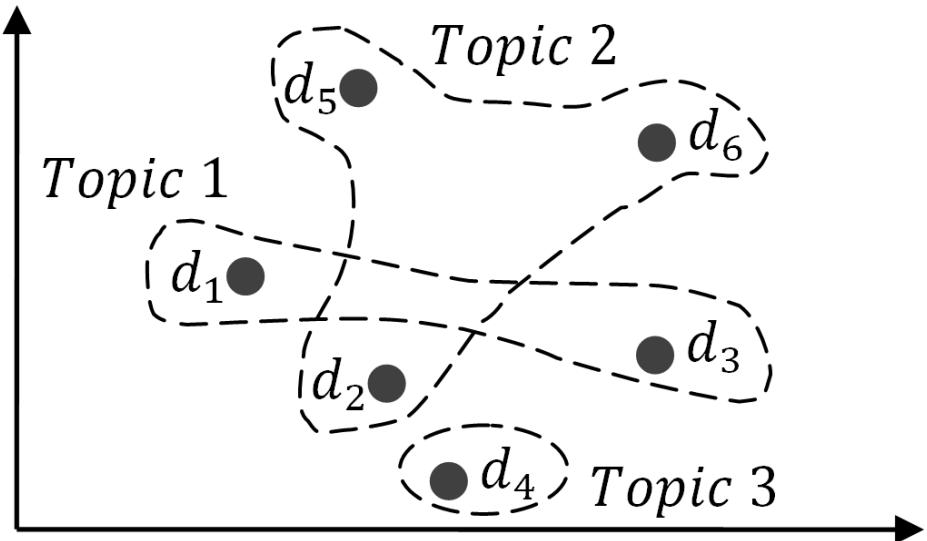
LDA Model 1

#Iterations₁; #topics₁; α_1 ; β_1



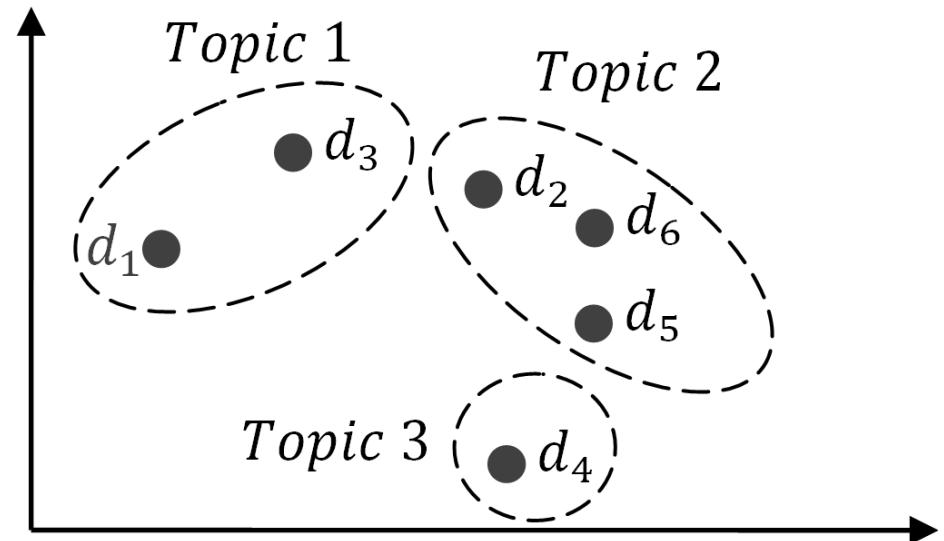
LDA Model 1

#Iterations₁; #topics₁; α_1 ; β_1



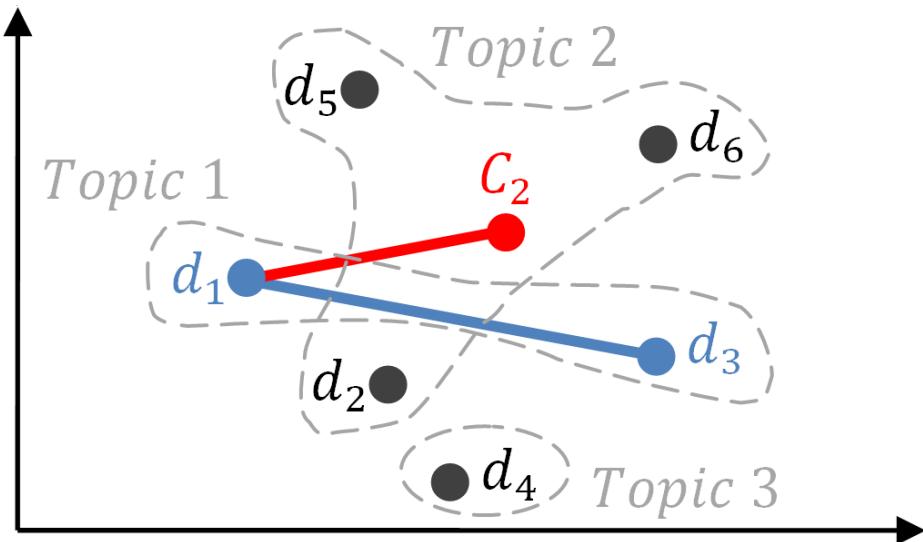
LDA Model 2

#Iterations₂; #topics₂; α_2 ; β_2



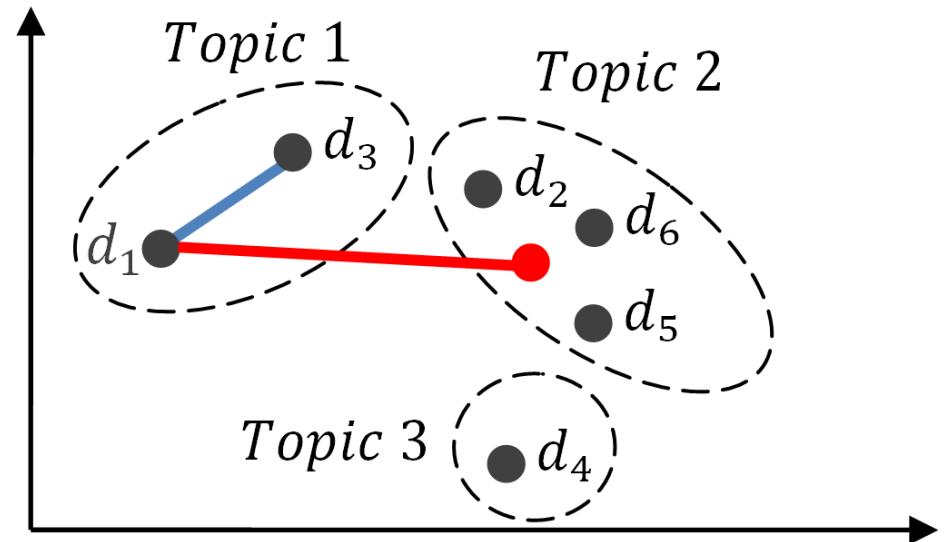
LDA Model 1

#Iterations₁; #topics₁; α_1 ; β_1



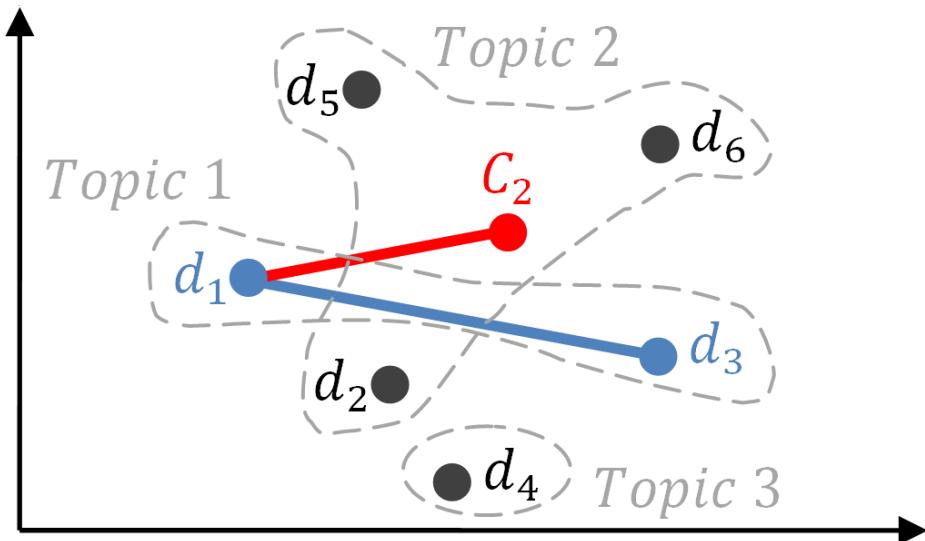
LDA Model 2

#Iterations₂; #topics₂; α_2 ; β_2



LDA Model 1

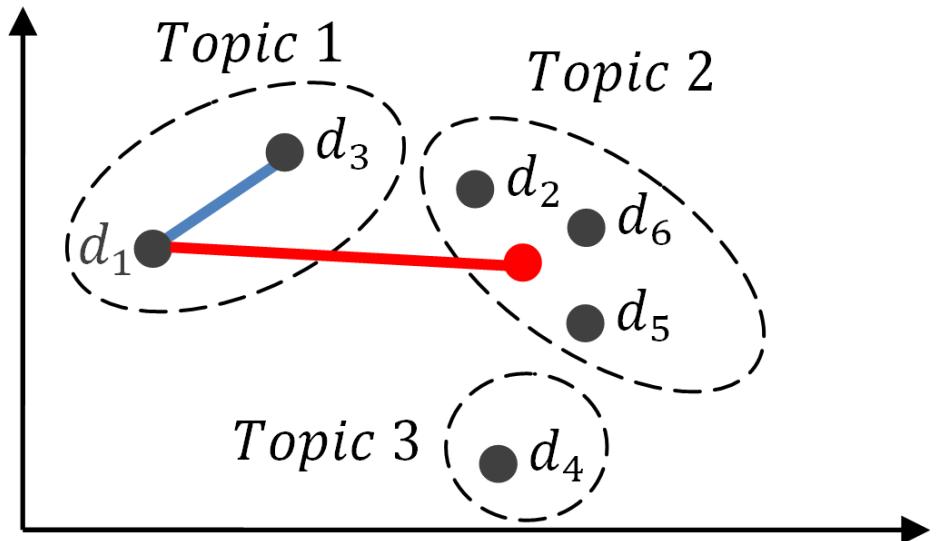
#Iterations₁; #topics₁; α_1 ; β_1



silhouette = 0.3

LDA Model 2

#Iterations₂; #topics₂; α_2 ; β_2



silhouette = 0.9

**Clusters more cohesive
Clusters well separated**

How to **evaluate** how “good” an
LDA configuration is?

How to **identify** the “good” LDA
parameter configurations?

How to identify the “good” LDA parameter configurations?

```
for numIter in [500, ...]  
    for numTopics in [5, ...]  
        for  $\alpha$  in [0.01, ...]  
            for  $\beta$  in [0.01, ...]  
                LDA[numIter , numTopics ,  $\alpha$ ,  $\beta$ ]
```

Exhaustive approach:
- Discretize search space & iterate?

How to identify the “good” LDA parameter configurations?

```
for numIter in [500, ...]  
    for numTopics in [5, ...]  
        for α in [0.01, ...]  
            for β in [0.01, ...]  
                LDA[numIter, numTopics, α, β]
```

Exhaustive approach:
- Discretize search space & iterate?

Too many possibilities

Use a Genetic Algorithm

What is a Genetic Algorithm (GA)?

- Stochastic search technique based on the process of natural evolution to identify *near-optimal solutions* to search problems

**Choose a random
population of LDA
parameters**

Choose a random population of LDA parameters

**First generation:
Population of random chromosomes**

| | iteration | topics | α | β |
|------------|-----------|--------|----------|---------|
| LDA Cfg. 1 | 510 | 74 | 0.34 | 2.5 |
| LDA Cfg. 2 | 725 | 128 | 1.28 | 0.4 |
| LDA Cfg. 3 | 814 | 97 | 0.43 | 0.9 |
| ... | ... | ... | ... | ... |
| LDA Cfg. n | 618 | 250 | 1.14 | 0.74 |

Choose a random population of LDA parameters

First generation:
Population of random chromosomes

Individual (chromosome):
Represents one possible LDA parameter configuration

| | iteration | topics | α | β |
|------------|-----------|--------|----------|---------|
| LDA Cfg. 1 | 510 | 74 | 0.34 | 2.5 |
| LDA Cfg. 2 | 725 | 128 | 1.28 | 0.4 |
| LDA Cfg. 3 | 814 | 97 | 0.43 | 0.9 |
| LDA Cfg. n | 618 | 250 | 1.14 | 0.74 |

Choose a random population of LDA parameters

First generation:
Population of random chromosomes

Individual (chromosome):
Represents one possible LDA parameter configuration

| | iteration | topics | α | β |
|------------|-----------|--------|----------|---------|
| LDA Cfg. 1 | 510 | 74 | 0.34 | 2.5 |
| LDA Cfg. 2 | 725 | 128 | 1.28 | 0.4 |
| LDA Cfg. 3 | 814 | 97 | 0.43 | 0.9 |
| LDA Cfg. n | 618 | 250 | 1.14 | 0.74 |

Gene:
LDA parameter

Choose a random population of LDA parameters

| | iteration | topics | α | β |
|------------|-----------|--------|----------|---------|
| LDA Cfg. 1 | 510 | 74 | 0.34 | 2.5 |
| LDA Cfg. 2 | 725 | 128 | 1.28 | 0.4 |
| LDA Cfg. 3 | 814 | 97 | 0.43 | 0.9 |
| ... | ... | ... | ... | ... |
| LDA Cfg. n | 618 | 250 | 1.14 | 0.74 |

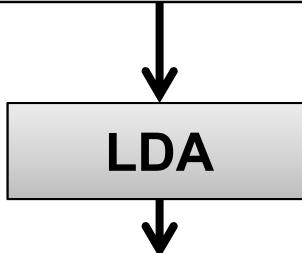
Choose a random population of LDA parameters



LDA

| | iteration | topics | α | β |
|-------------------|-----------|--------|----------|---------|
| LDA Cfg. 1 | 510 | 74 | 0.34 | 2.5 |
| LDA Cfg. 2 | 725 | 128 | 1.28 | 0.4 |
| LDA Cfg. 3 | 814 | 97 | 0.43 | 0.9 |
| ... | ... | ... | ... | ... |
| LDA Cfg. n | 618 | 250 | 1.14 | 0.74 |

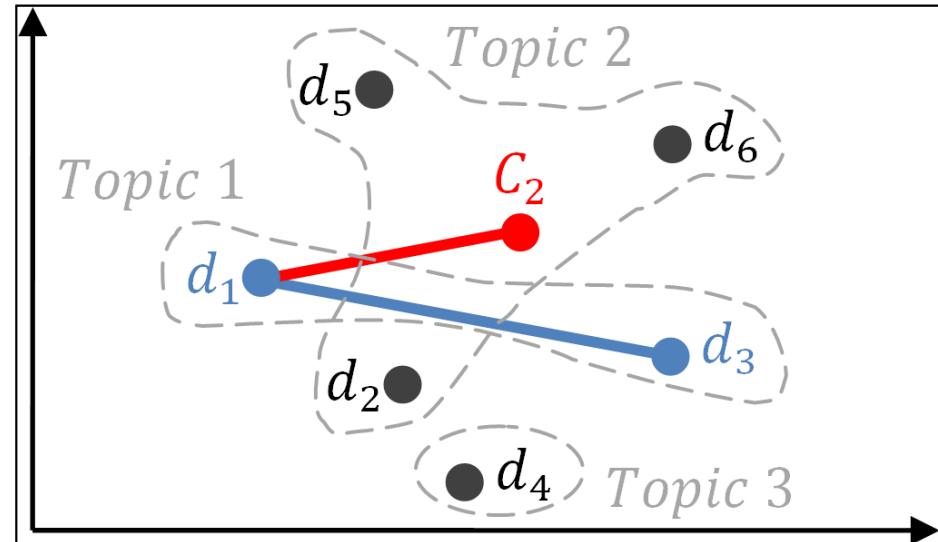
Choose a random population of LDA parameters



| | iteration | topics | α | β |
|------------|-----------|--------|----------|---------|
| LDA Cfg. 1 | 510 | 74 | 0.34 | 2.5 |
| LDA Cfg. 2 | 725 | 128 | 1.28 | 0.4 |
| LDA Cfg. 3 | 814 | 97 | 0.43 | 0.9 |
| ... | ... | ... | ... | ... |
| LDA Cfg. n | 618 | 250 | 1.14 | 0.74 |

Determine **fitness** of each *chromosome* (individual)

Fitness = silhouette coefficient



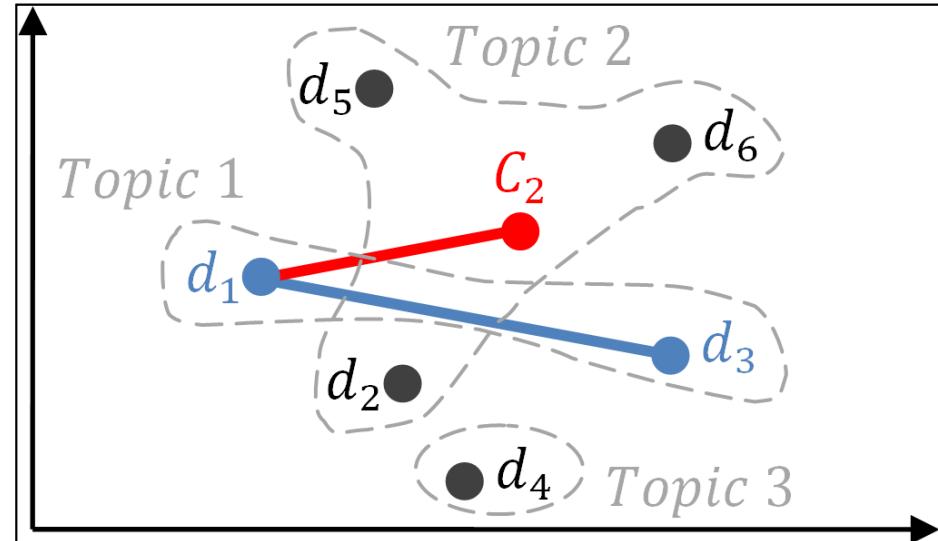
Choose a random population of LDA parameters

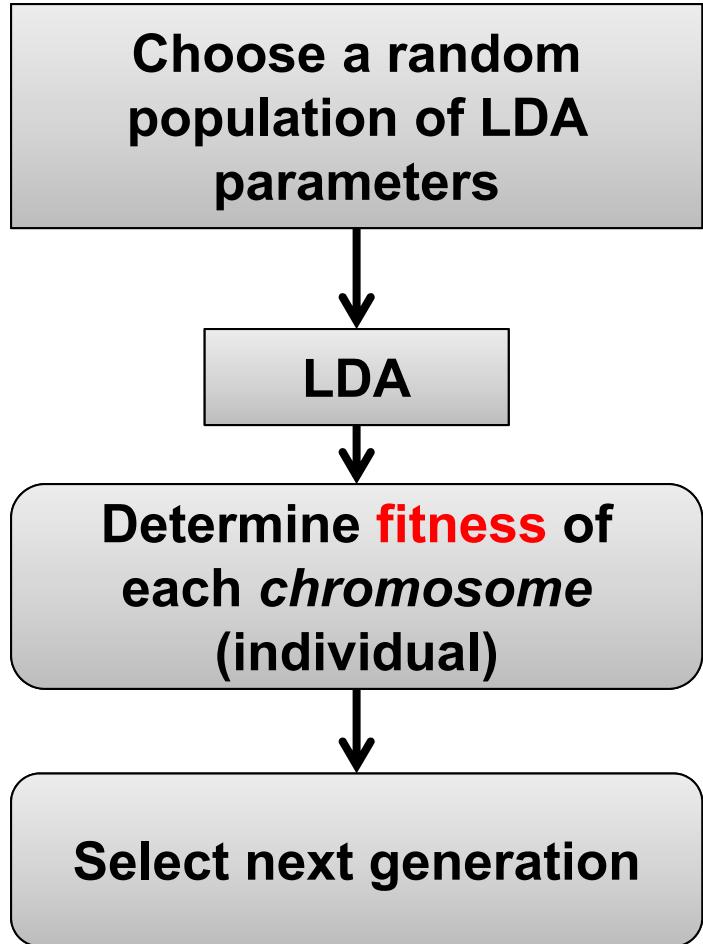
LDA

Determine **fitness** of each *chromosome* (individual)

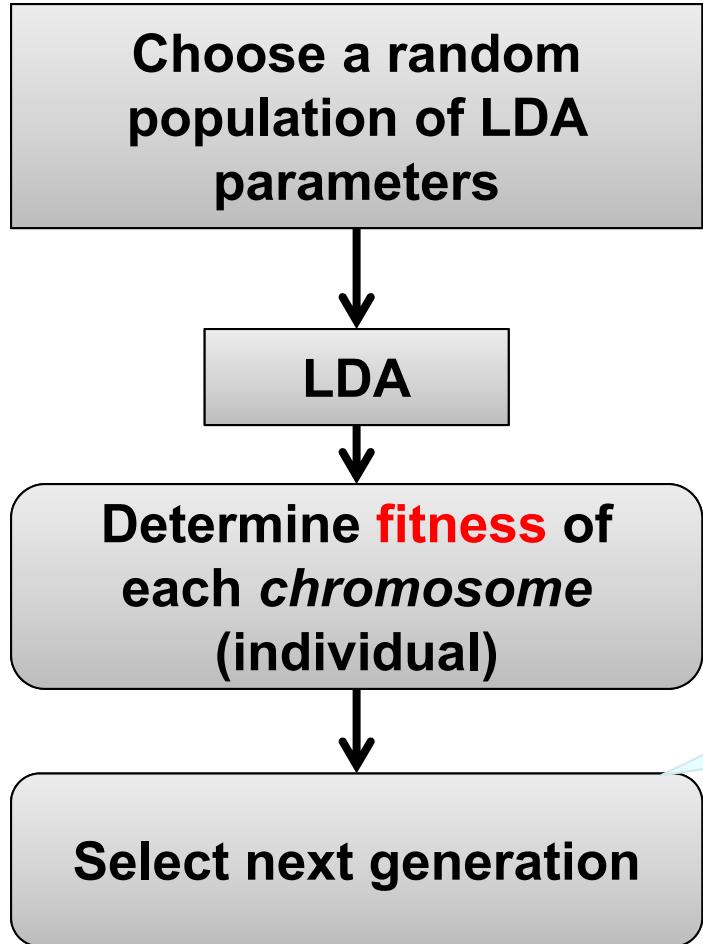
| | iteration | topics | α | β | Fitness |
|------------|-----------|--------|----------|---------|---------|
| LDA Cfg. 1 | 510 | 74 | 0.34 | 2.5 | 0.2 |
| LDA Cfg. 2 | 725 | 128 | 1.28 | 0.4 | 0.4 |
| LDA Cfg. 3 | 814 | 97 | 0.43 | 0.9 | 0.35 |
| ... | ... | ... | ... | ... | ... |
| LDA Cfg. n | 618 | 250 | 1.14 | 0.74 | 0.1 |

Fitness = silhouette coefficient



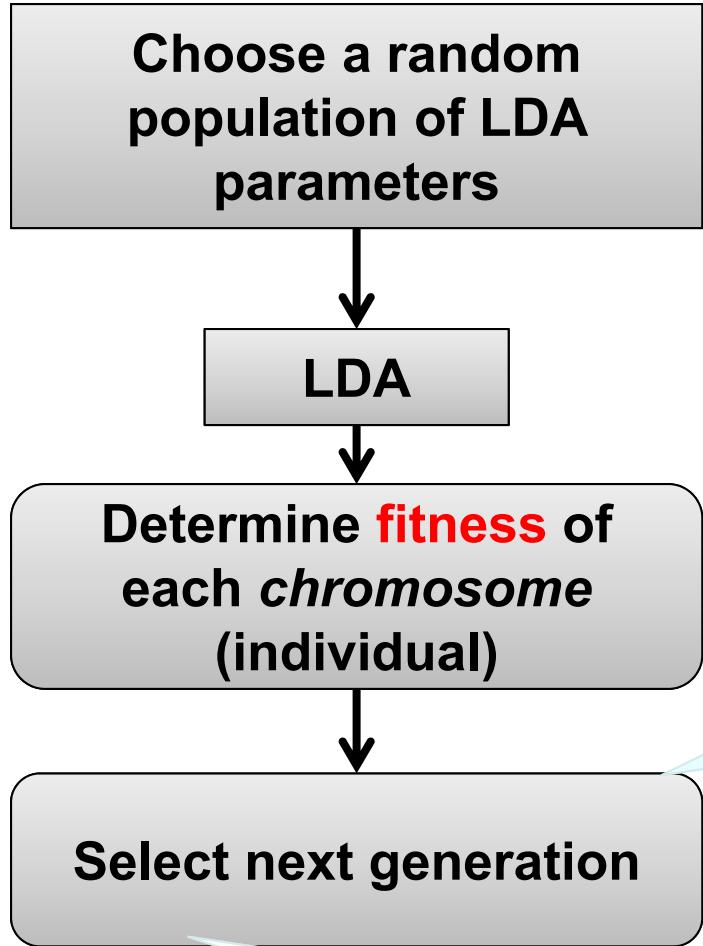


| | iteration | topics | α | β | Fitness |
|------------|-----------|--------|----------|---------|---------|
| LDA Cfg. 1 | 510 | 74 | 0.34 | 2.5 | 0.2 |
| LDA Cfg. 2 | 725 | 128 | 1.28 | 0.4 | 0.4 |
| LDA Cfg. 3 | 814 | 97 | 0.43 | 0.9 | 0.35 |
| ... | ... | ... | ... | ... | ... |
| LDA Cfg. n | 618 | 250 | 1.14 | 0.74 | 0.1 |



| | iteration | topics | α | β | Fitness |
|------------|-----------|--------|----------|---------|---------|
| LDA Cfg. 1 | 510 | 74 | 0.34 | 2.5 | 0.2 |
| LDA Cfg. 2 | 725 | 128 | 1.28 | 0.4 | 0.4 |
| LDA Cfg. 3 | 814 | 97 | 0.43 | 0.9 | 0.35 |
| ... | ... | ... | ... | ... | ... |
| LDA Cfg. n | 618 | 250 | 1.14 | 0.7 | 0.1 |

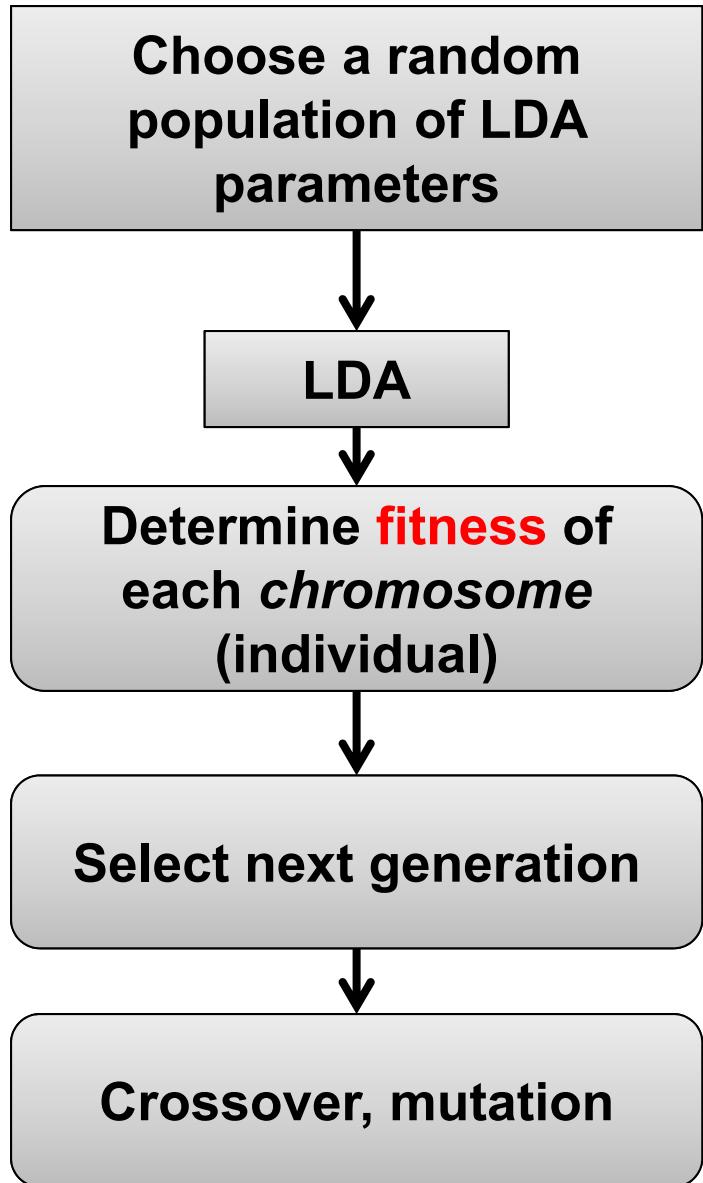
Elitism:
best (n=2) configurations will survive for next generation

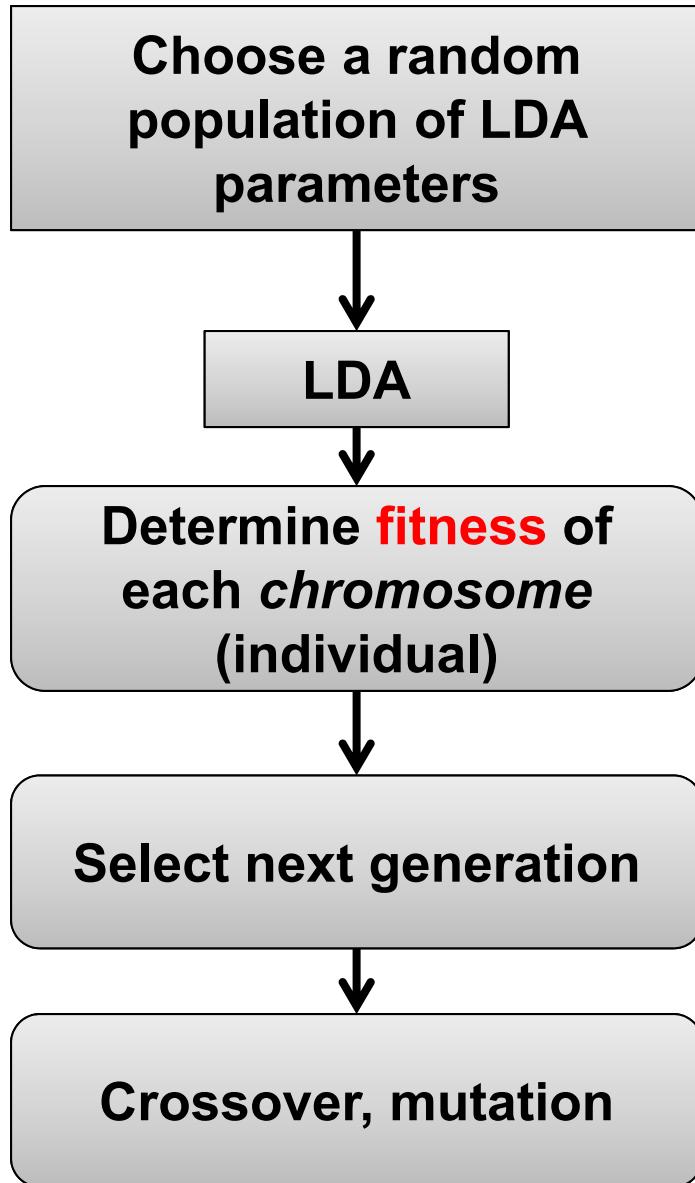


| | iteration | topics | α | β | Fitness |
|------------|-----------|--------|----------|---------|---------|
| LDA Cfg. 1 | 510 | 74 | 0.34 | 2.5 | 0.2 |
| LDA Cfg. 2 | 725 | 128 | 1.28 | 0.4 | 0.4 |
| LDA Cfg. 3 | 814 | 97 | 0.43 | 0.9 | 0.35 |
| ... | ... | ... | ... | ... | ... |
| LDA Cfg. n | 618 | 250 | 1.14 | 0.7 | 0.1 |

Elitism:
best (n=2) configurations will survive for next generation

Roulette selection:
Chance of chromosomes to contribute to next generation is proportional to their **fitness**

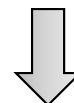




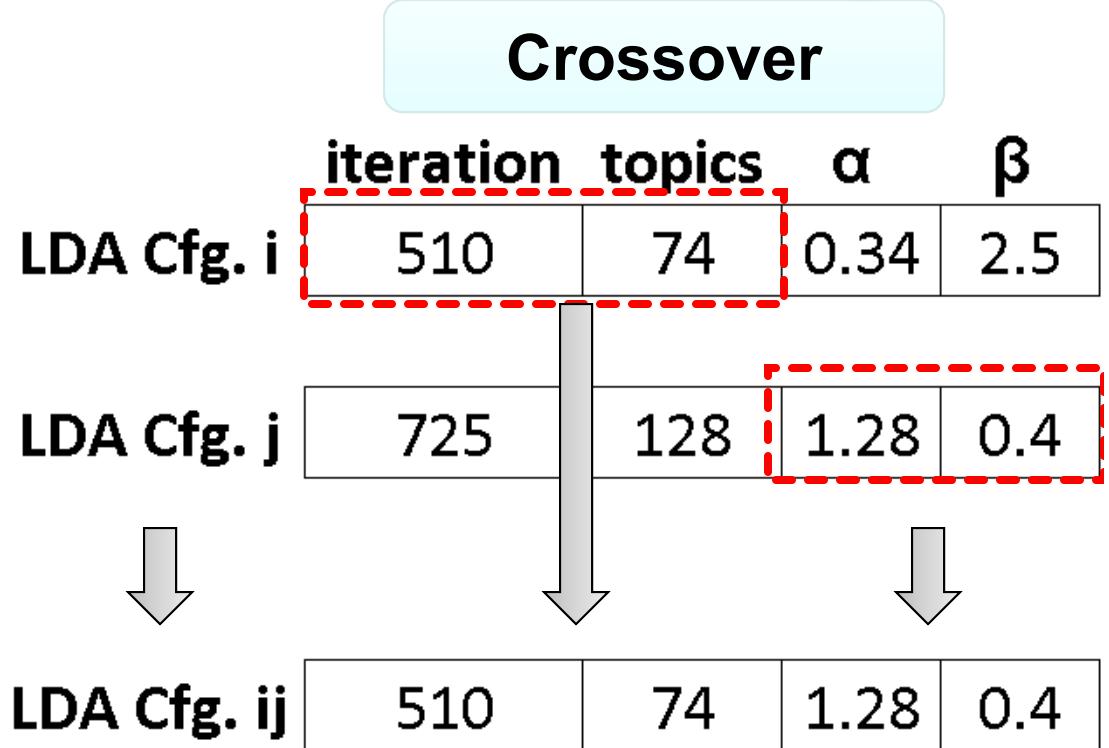
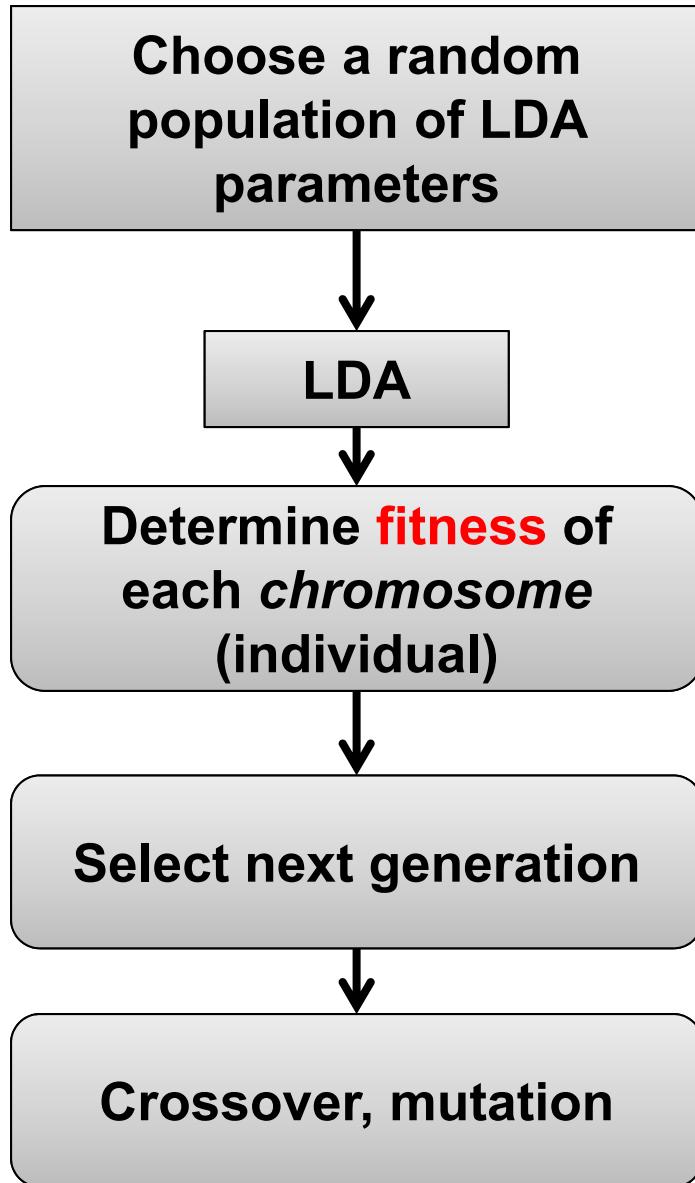
Crossover

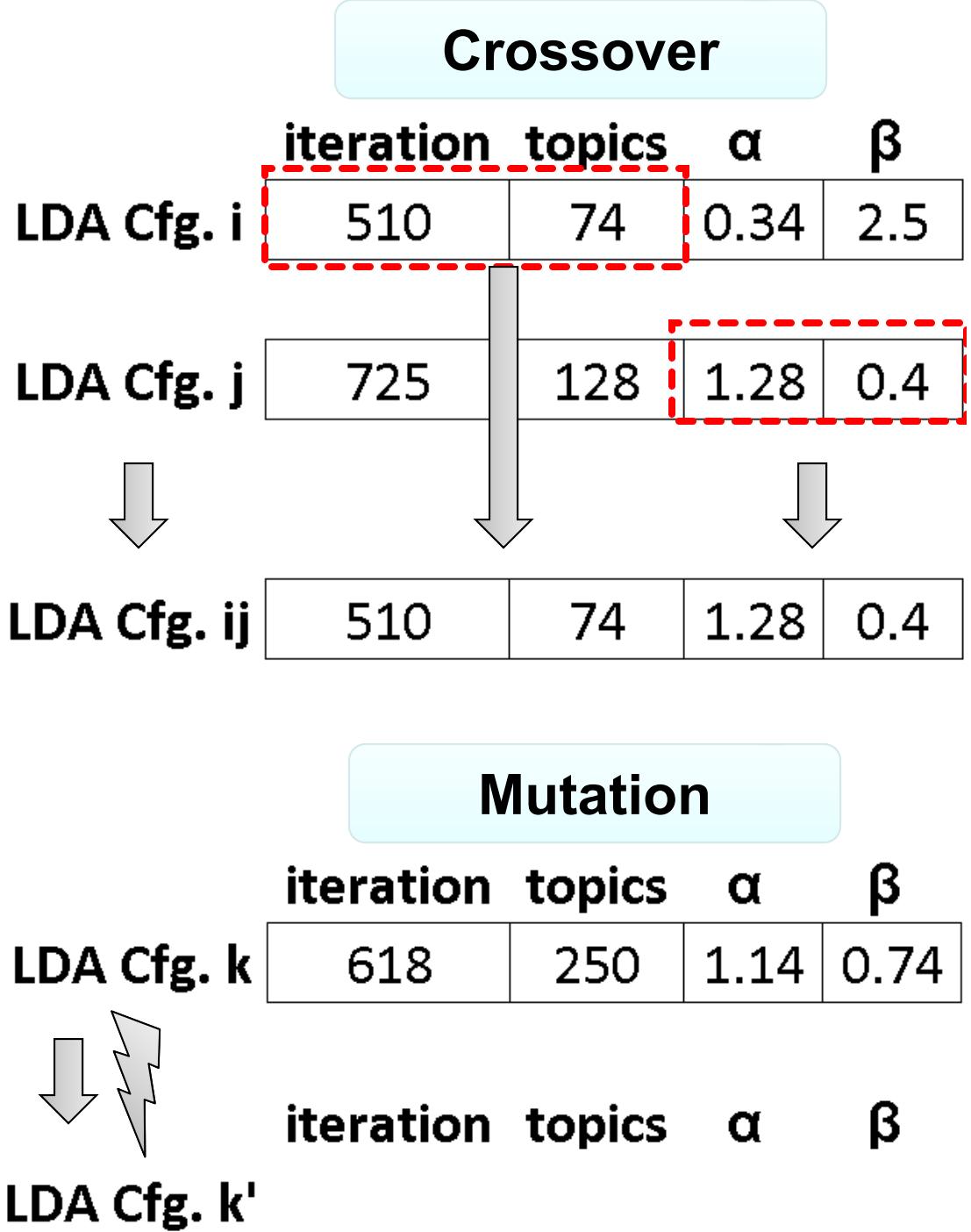
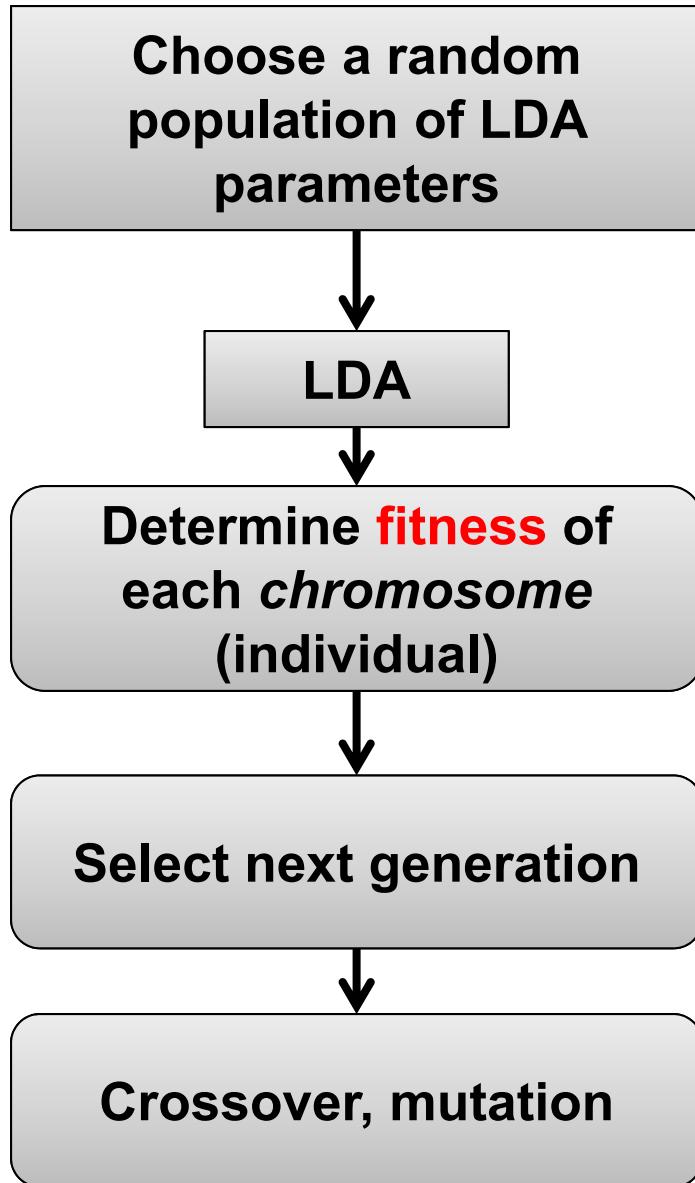
| | iteration | topics | α | β |
|------------|-----------|--------|----------|---------|
| LDA Cfg. i | 510 | 74 | 0.34 | 2.5 |

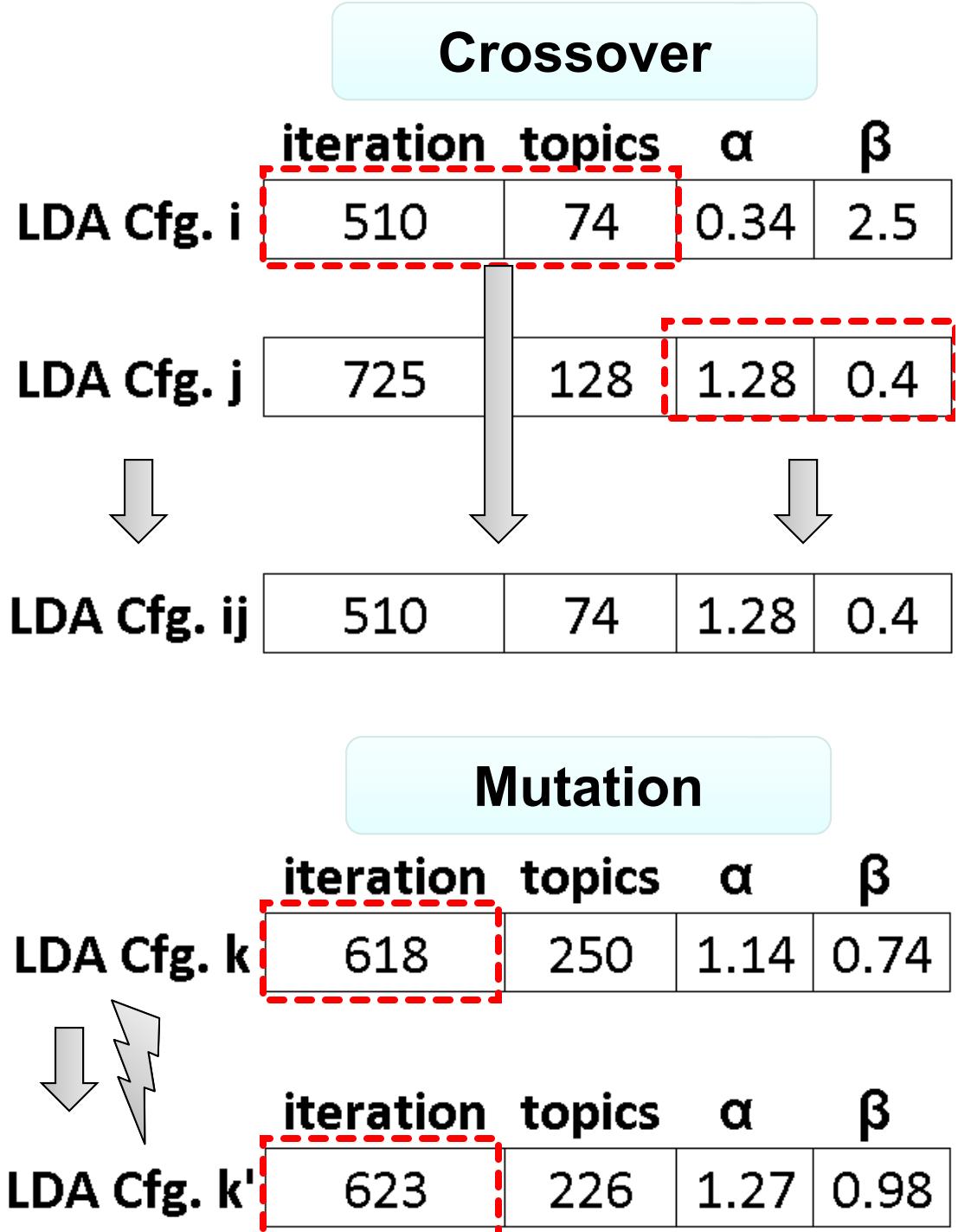
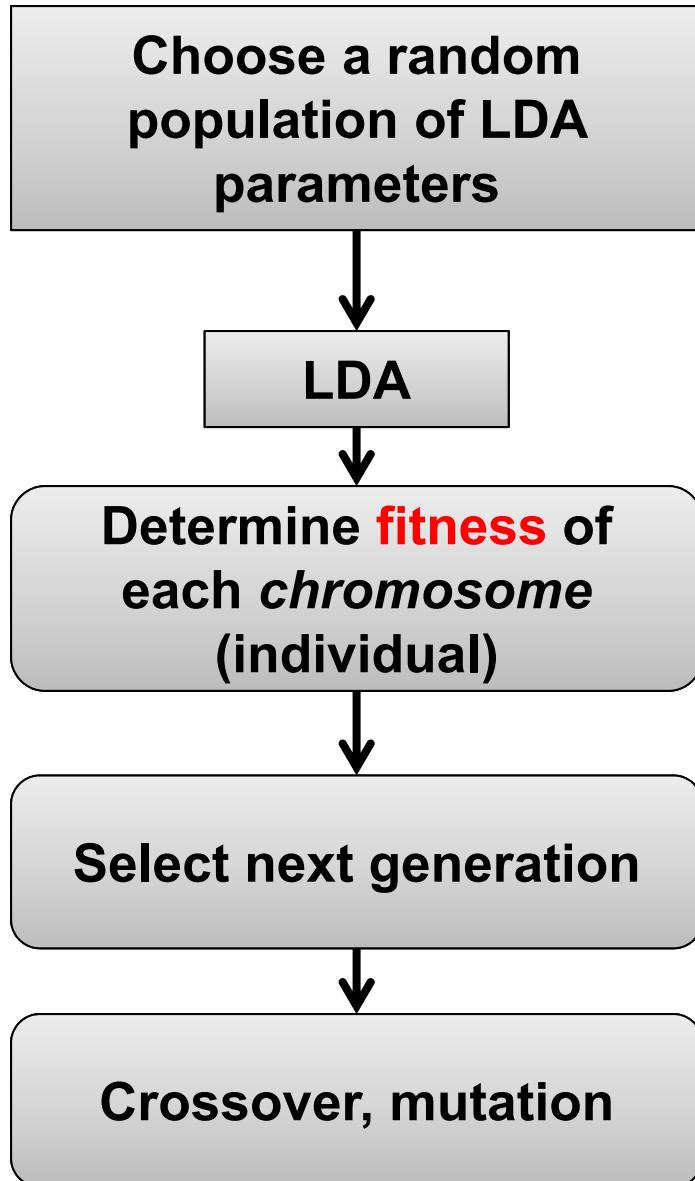
| | | | | |
|------------|-----|-----|------|-----|
| LDA Cfg. j | 725 | 128 | 1.28 | 0.4 |
|------------|-----|-----|------|-----|

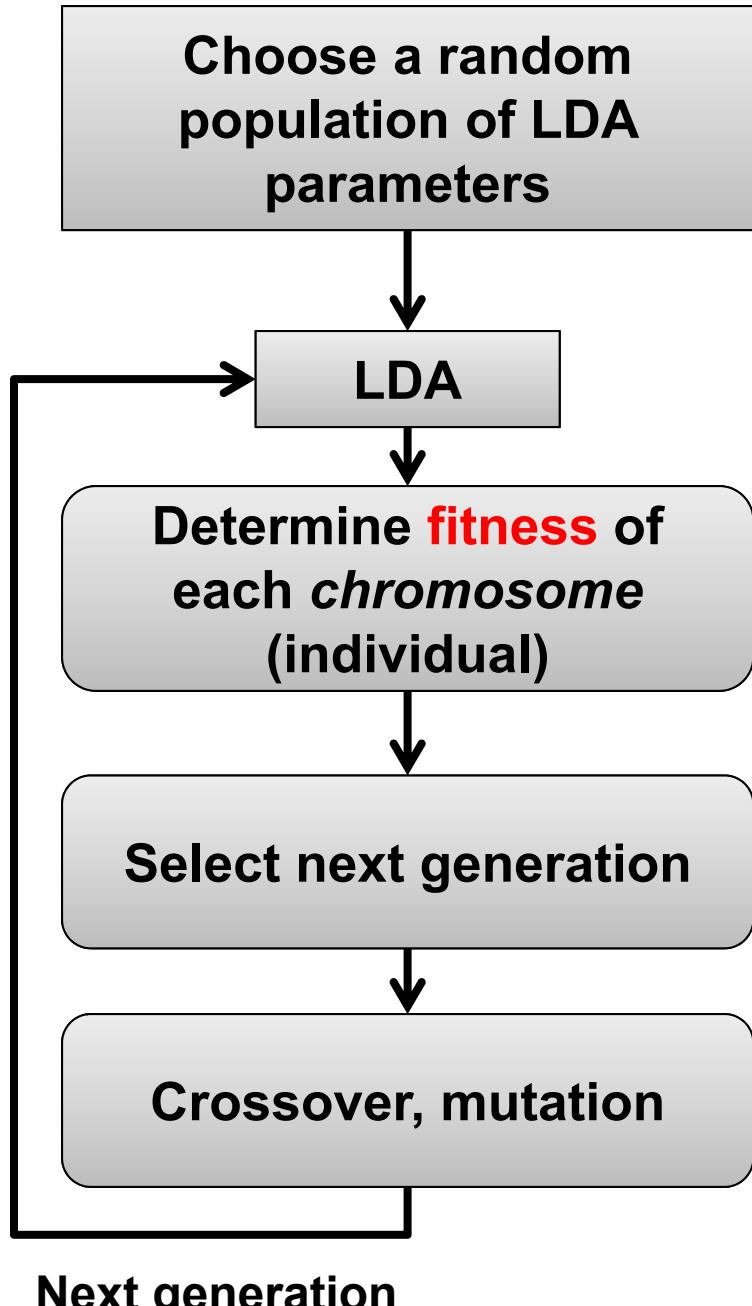


LDA Cfg. ij

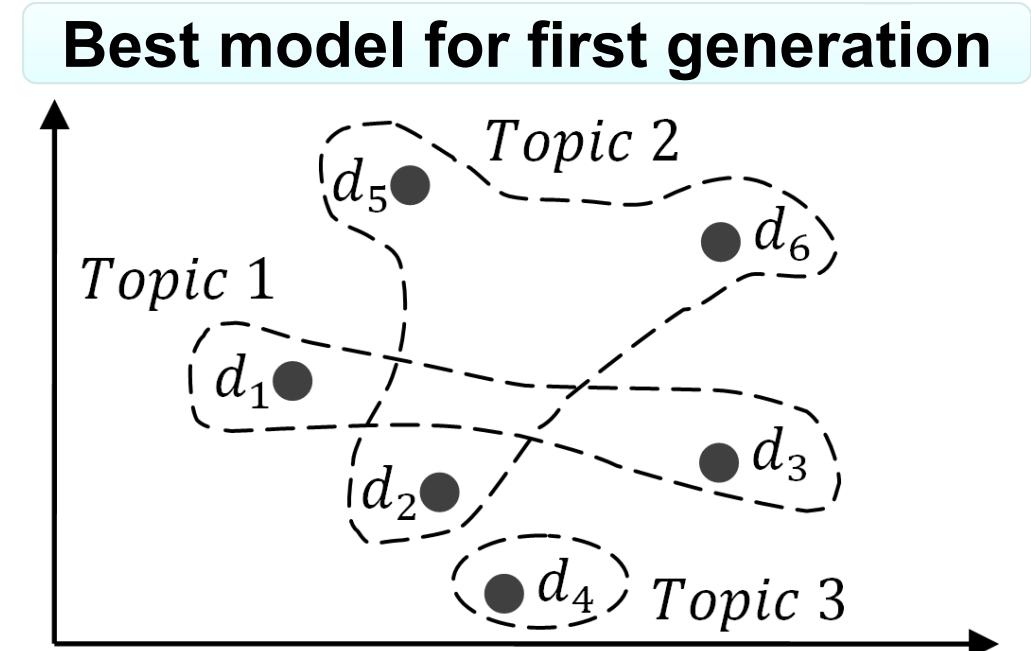
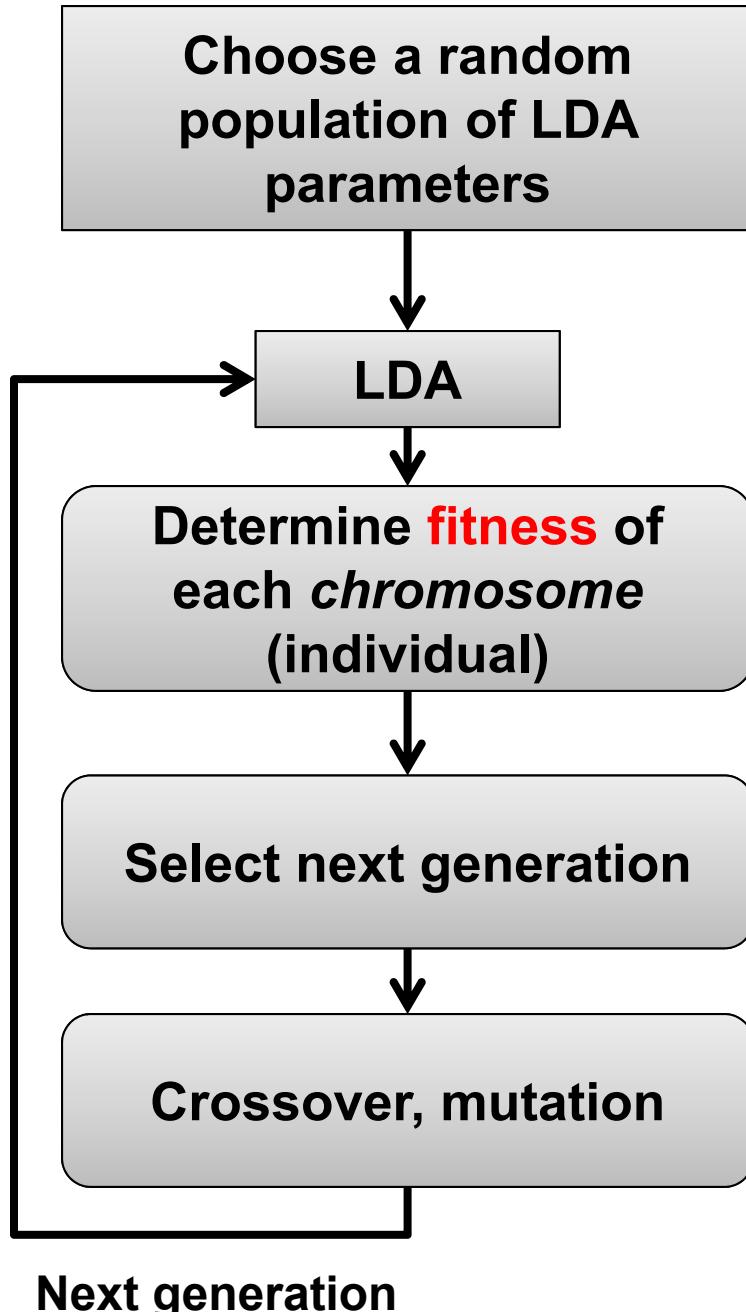


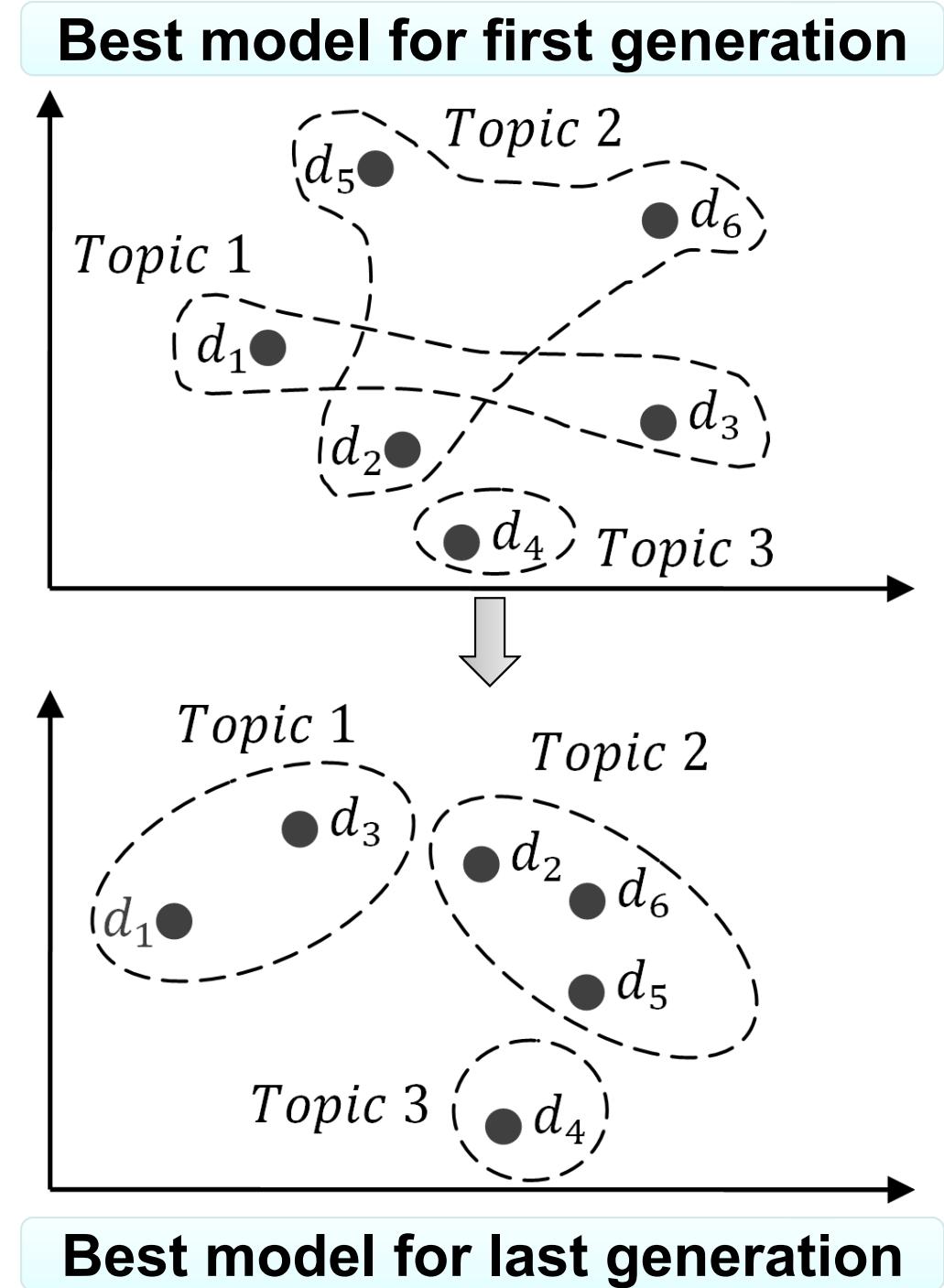
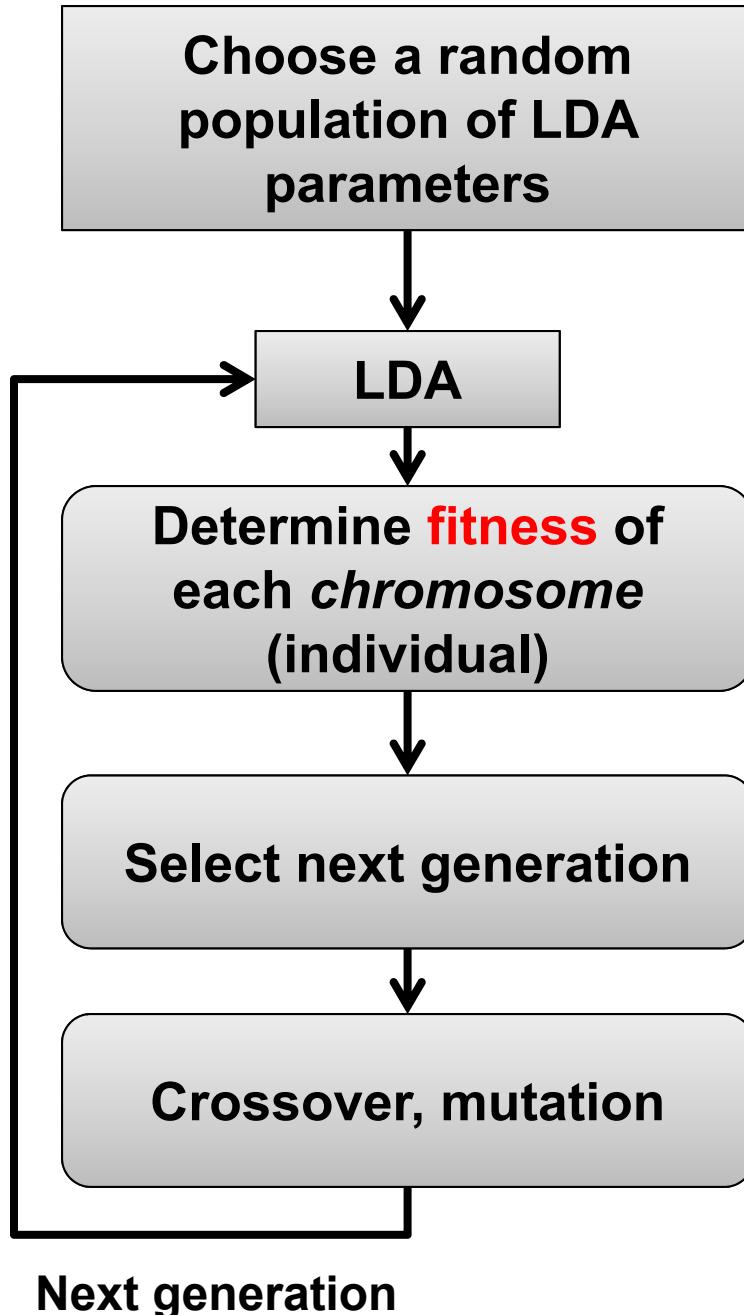


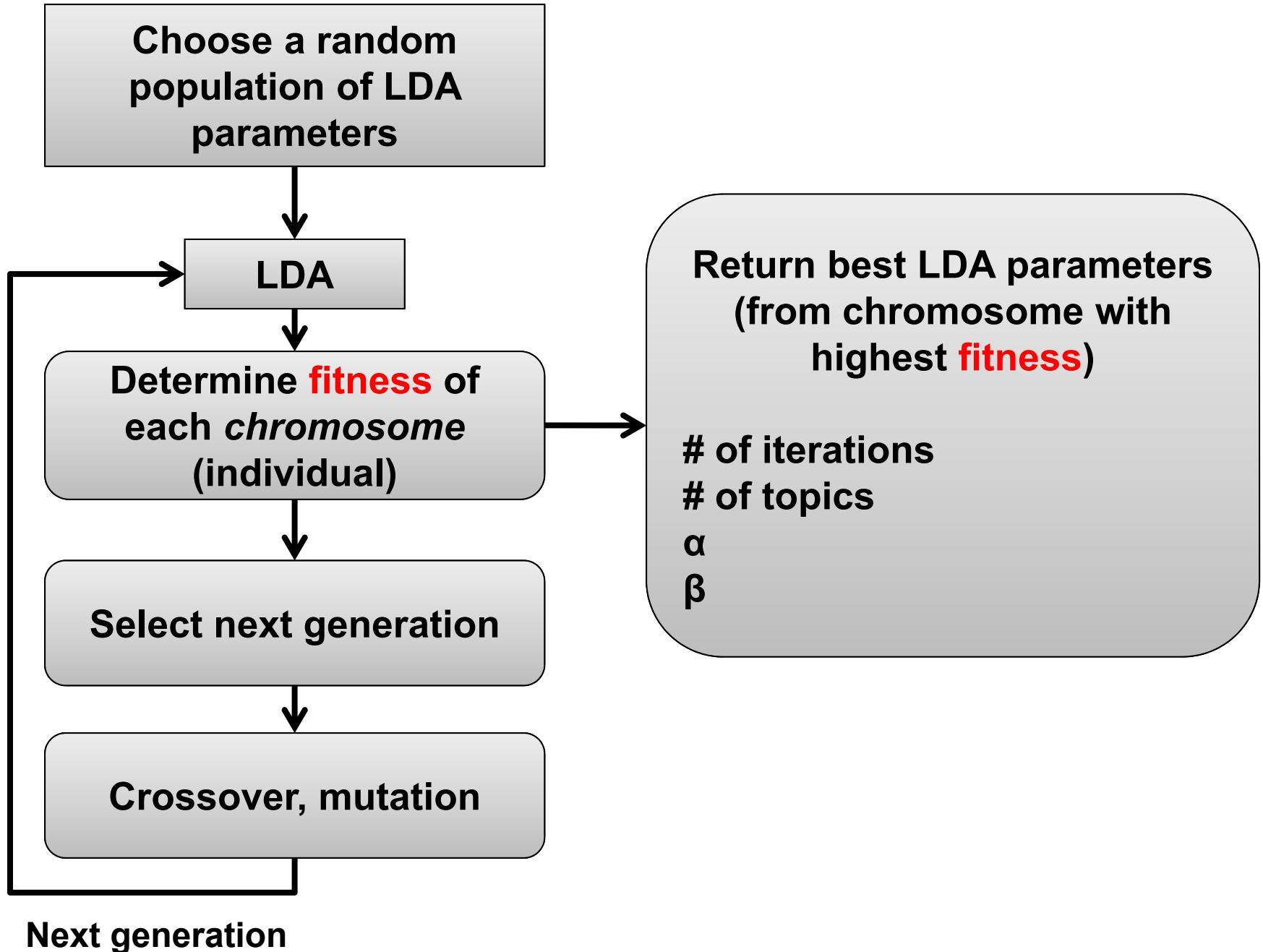




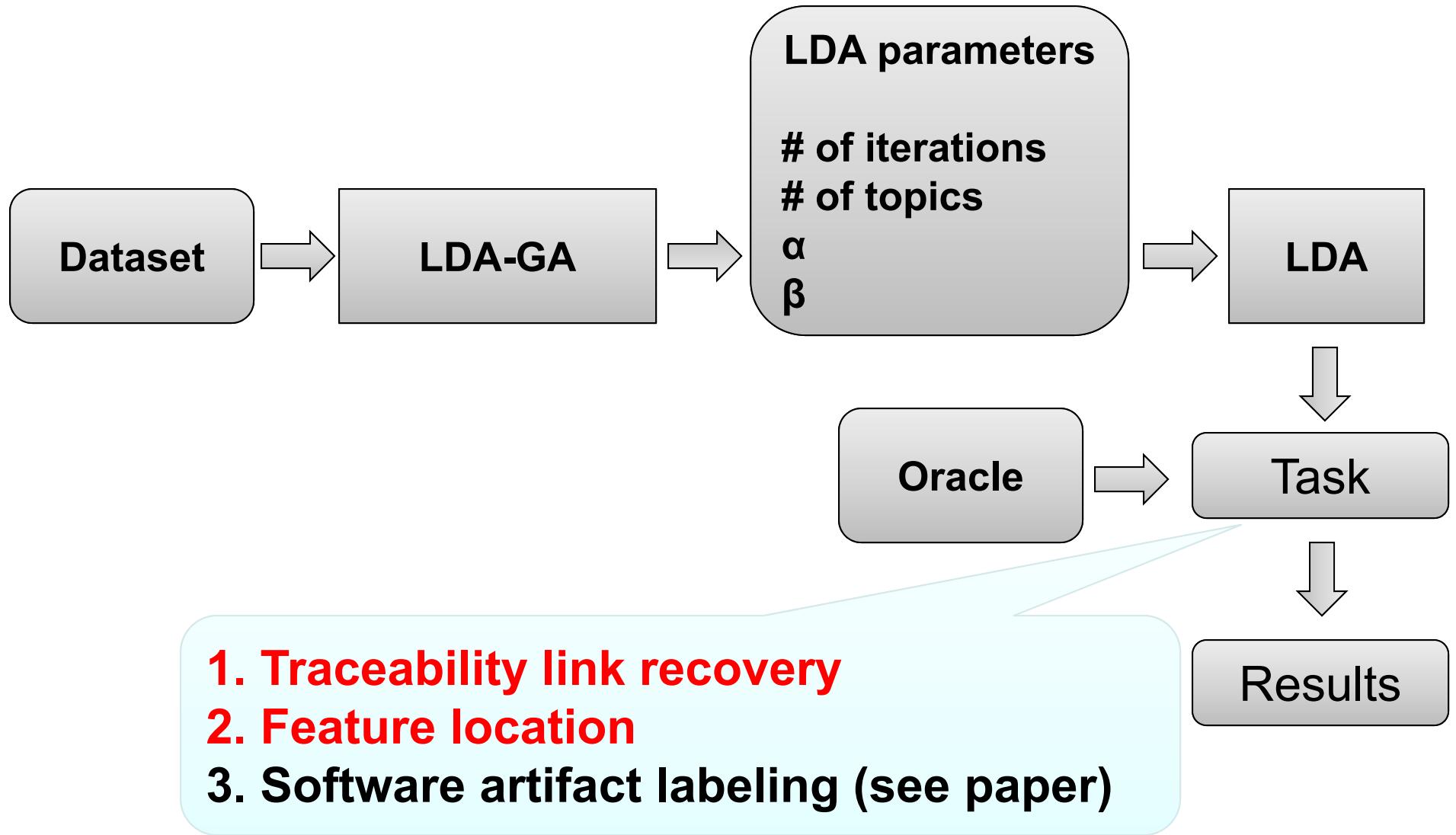
Next generation







Evaluation...



Evaluation: Traceability Link Recovery

- Recover links between *use cases* and *code classes*

| System | Size | # use cases | # code classes | # correct links |
|------------|--------|-------------|----------------|-----------------|
| EasyClinic | 20KLOC | 30 | 47 | 93 |
| eTour | 45KLOC | 58 | 174 | 366 |

Combinatorial:

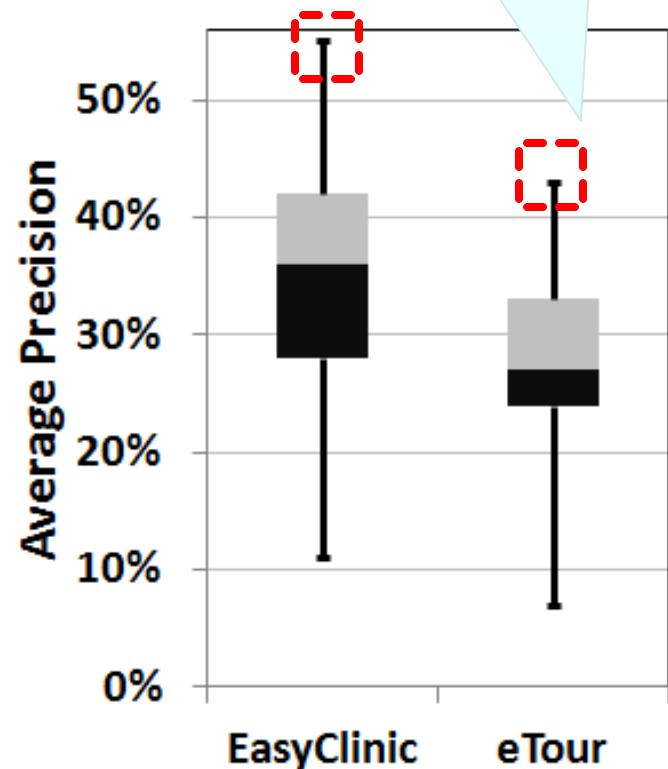
```
for numIter in [500, ...]  
    for numTopics in [5, ...]  
        for  $\alpha$  in [0.01, ...]  
            for  $\beta$  in [0.01, ...]  
                LDA[numIter , numTopics ,  $\alpha$ ,  $\beta$ ]
```

Choose LDA parameters
with best average
precision using an *oracle*

Combinatorial:

```
for numIter in [500, ...]  
    for numTopics in [5, ...]  
        for  $\alpha$  in [0.01, ...]  
            for  $\beta$  in [0.01, ...]  
                LDA[numIter , numTopics ,  $\alpha$  ,  $\beta$ ]
```

Choose LDA parameters
with best average
precision using an *oracle*



Combinatorial:

```
for numIter in [500, ...]  
    for numTopics in [5, ...]  
        for  $\alpha$  in [0.01, ...]  
            for  $\beta$  in [0.01, ...]  
                LDA[numIter , numTopics ,  $\alpha$ ,  $\beta$ ]
```

Choose LDA parameters
with best average
precision using an *oracle*

LDA-GA:

run LDA-GA 30 times (to
account for randomness)

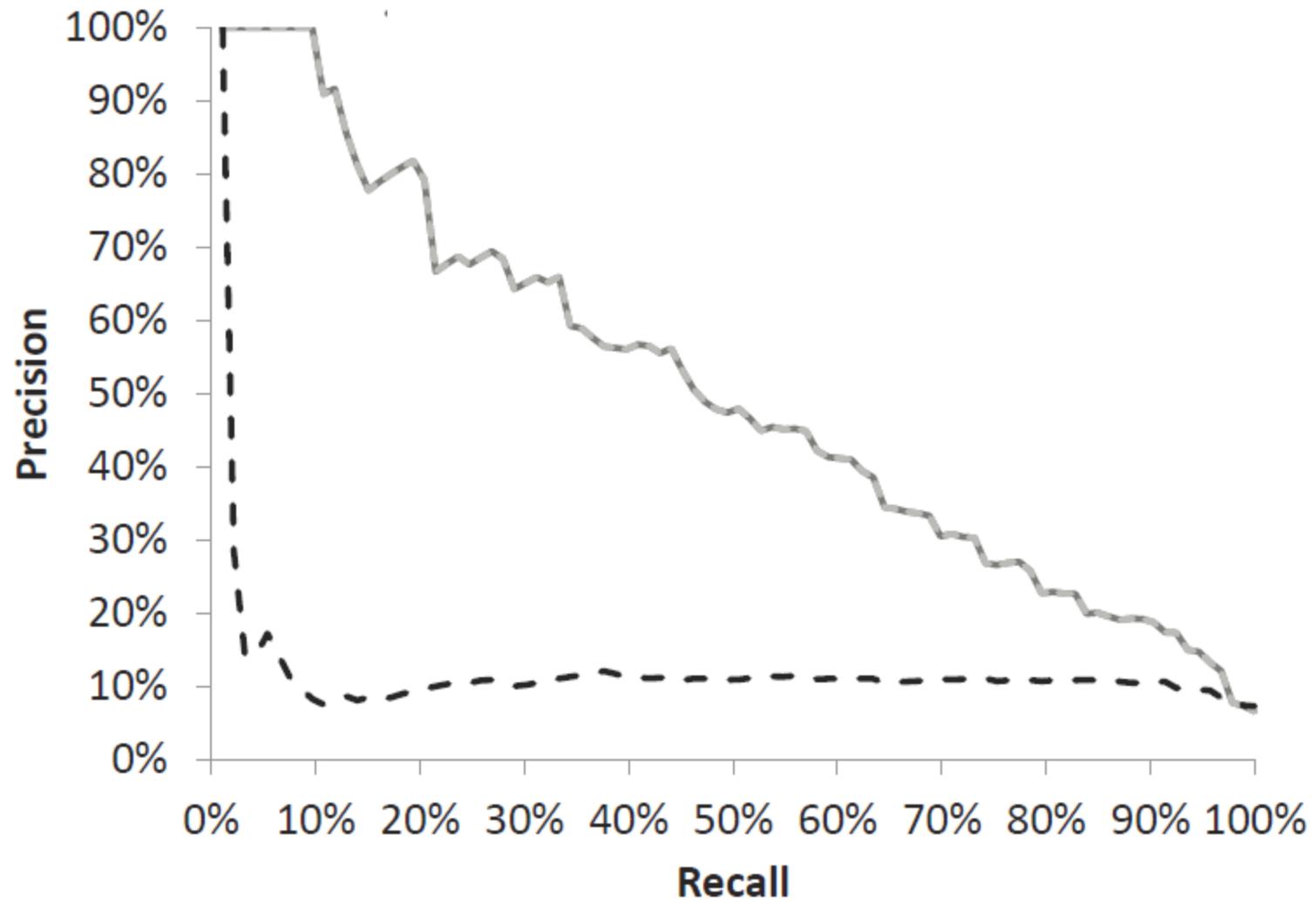
Choose LDA parameters
corresponding to median
fitness over 30 runs

Baseline:

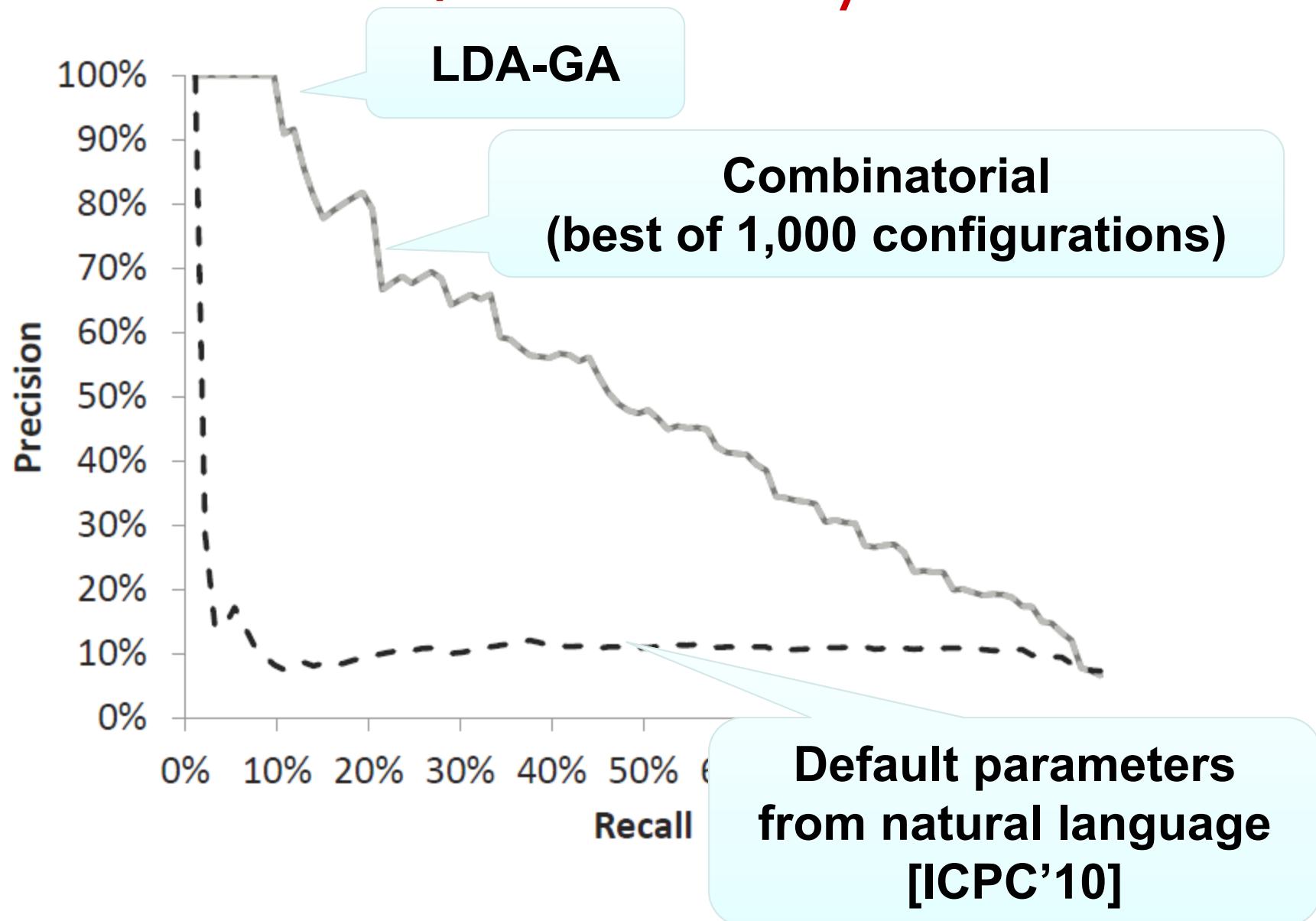
[ICPC'10]

Use default LDA
parameters from natural
language

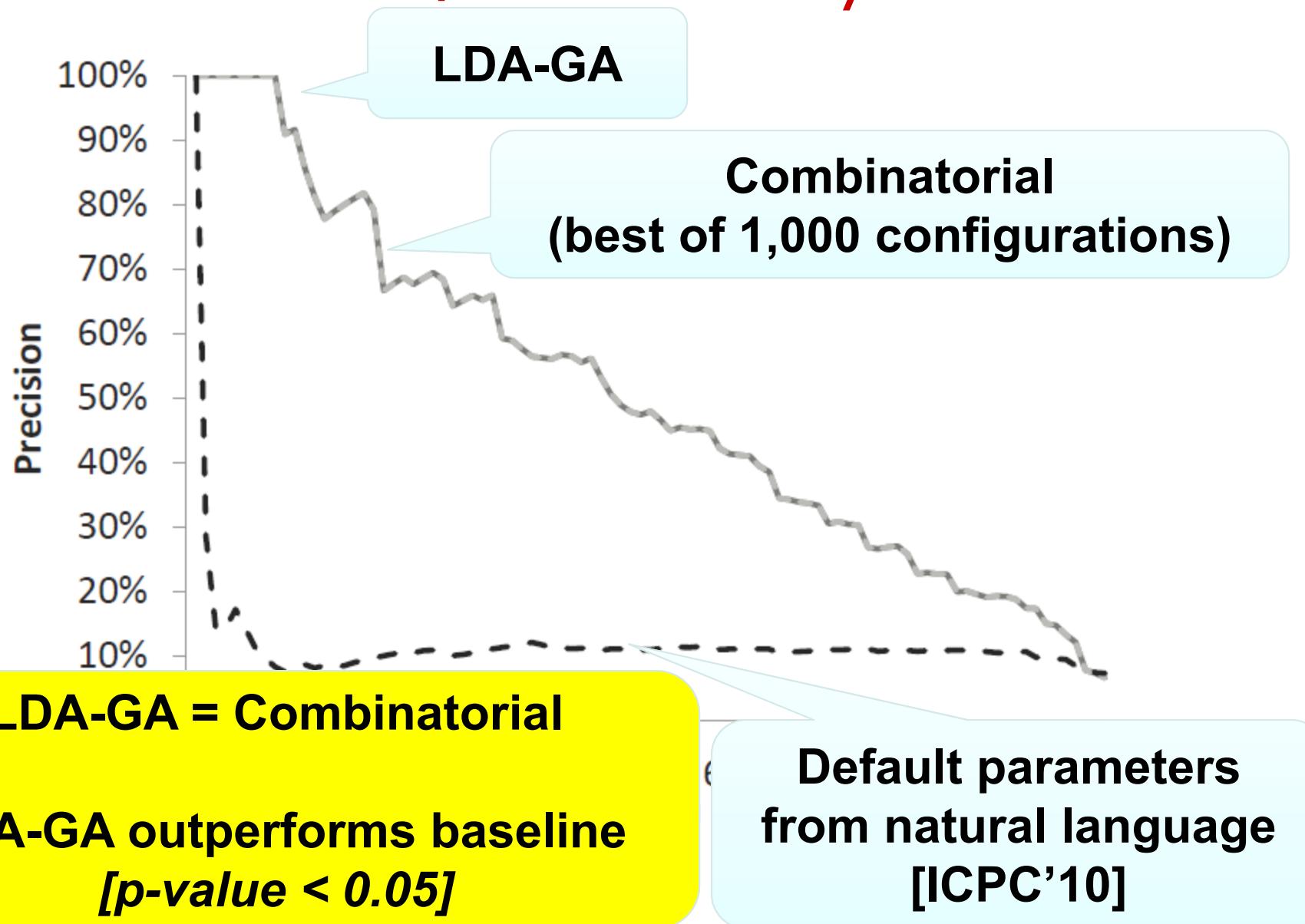
Precision/Recall EasyClinic



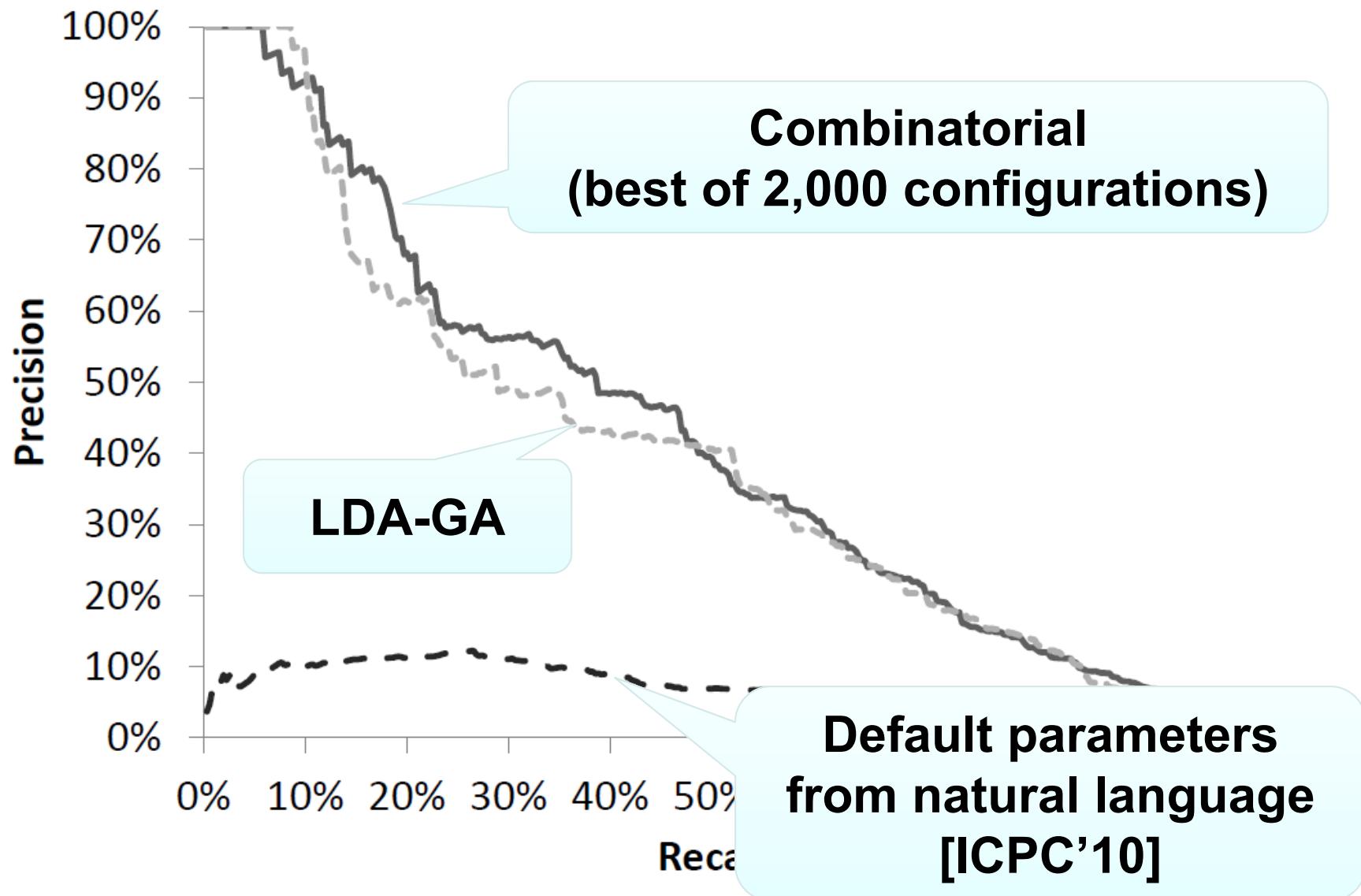
Precision/Recall EasyClinic



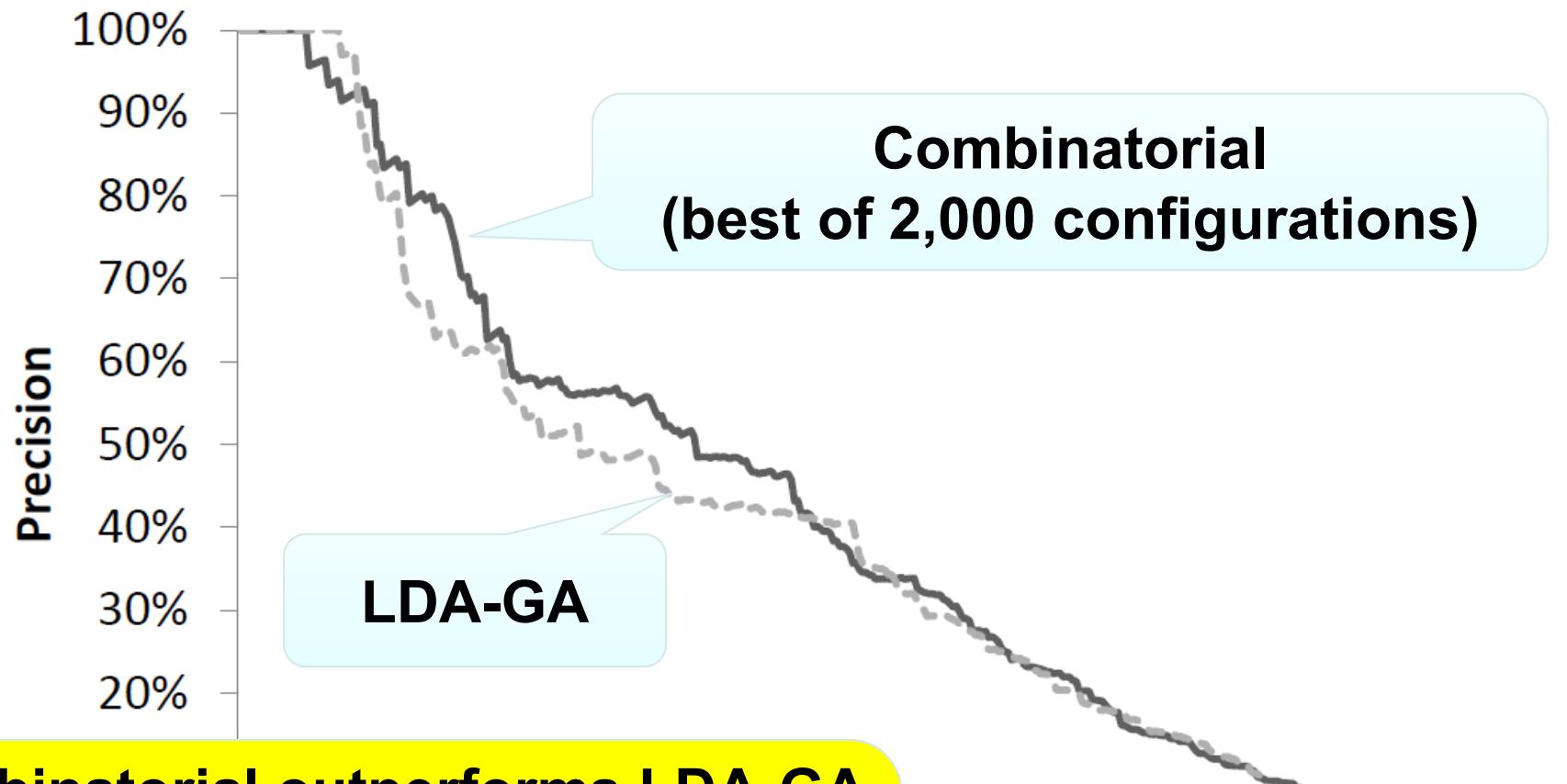
Precision/Recall EasyClinic



Precision/Recall eTour



Precision/Recall eTour



Combinatorial outperforms LDA-GA
[$p\text{-value} < 0.05$]

LDA-GA outperforms baseline

Default parameters
from natural language
[ICPC'10]

Evaluation: Feature location

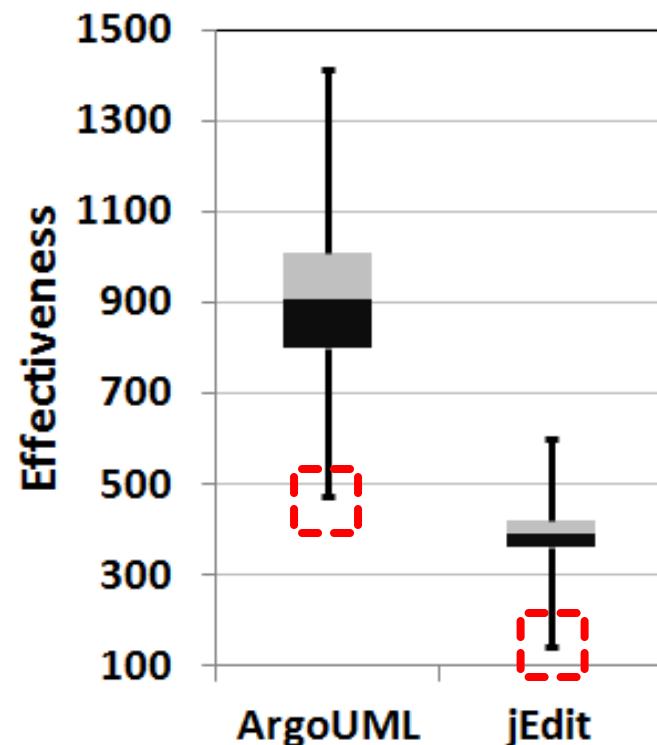
- Identify methods related to a maintenance task (e.g., bug, feature)

| System | Size | # features | # methods |
|---------------|-------------|-------------------|------------------|
| jEdit | 104KLOC | 150 | 6,413 |
| ArgoUML | 149KLOC | 91 | 11,000 |

Combinatorial:

```
for numIter in [500, ...]  
    for numTopics in [5, ...]  
        for  $\alpha$  in [0.01, ...]  
            for  $\beta$  in [0.01, ...]  
                LDA[numIter , numTopics ,  $\alpha$  ,  $\beta$ ]
```

Choose LDA parameters with
best average effectiveness
using an *oracle*



Combinatorial:

```
for numIter in [500, ...]  
    for numTopics in [5, ...]  
        for  $\alpha$  in [0.01, ...]  
            for  $\beta$  in [0.01, ...]  
                LDA[numIter , numTopics ,  $\alpha$ ,  $\beta$ ]
```

Choose LDA parameters with
best average effectiveness
using an *oracle*

LDA-GA:

run LDA-GA 30 times (to
account for randomness)

Choose LDA parameters
corresponding to median
fitness over 30 runs

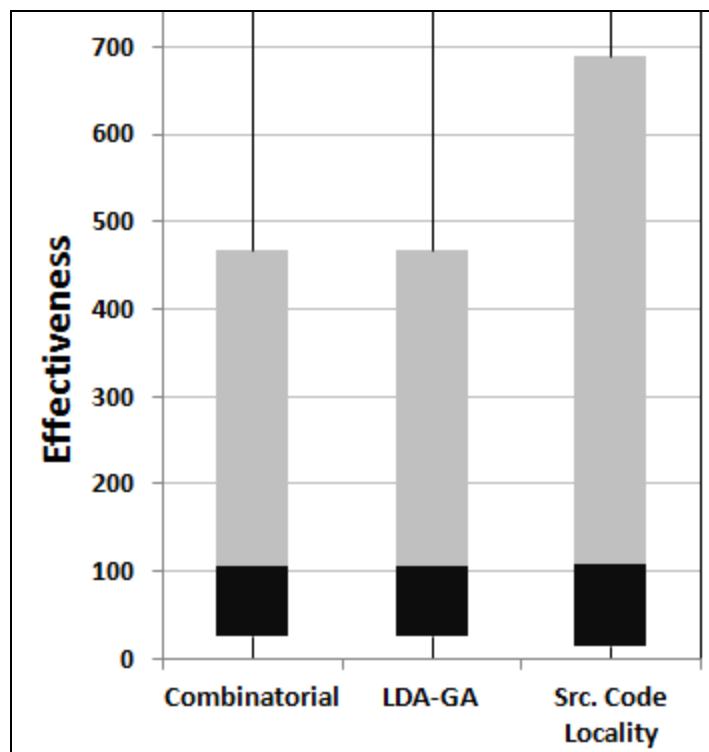
Baseline:

[SCAM'10]

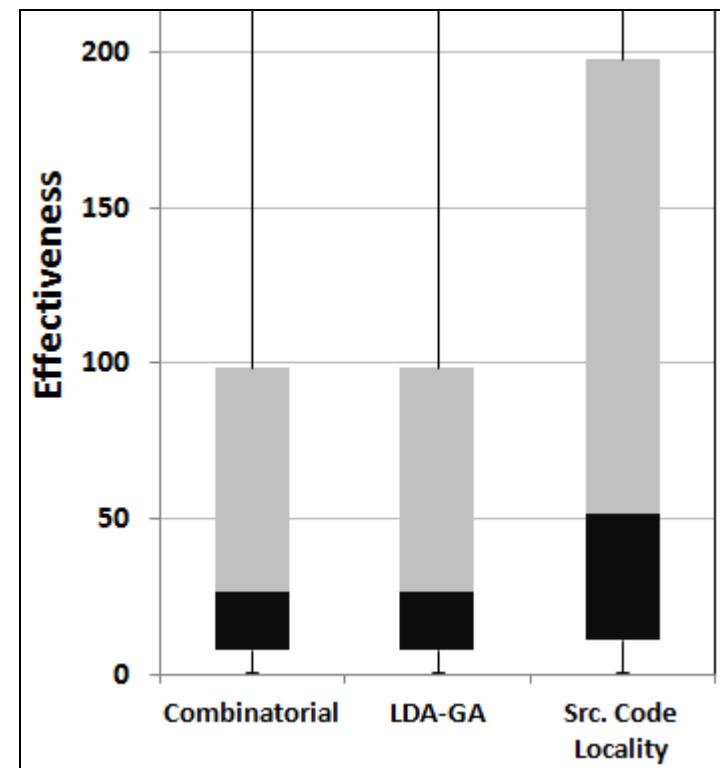
Use LDA parameters from
source locality heuristic

Effectiveness measure

ArgoUML

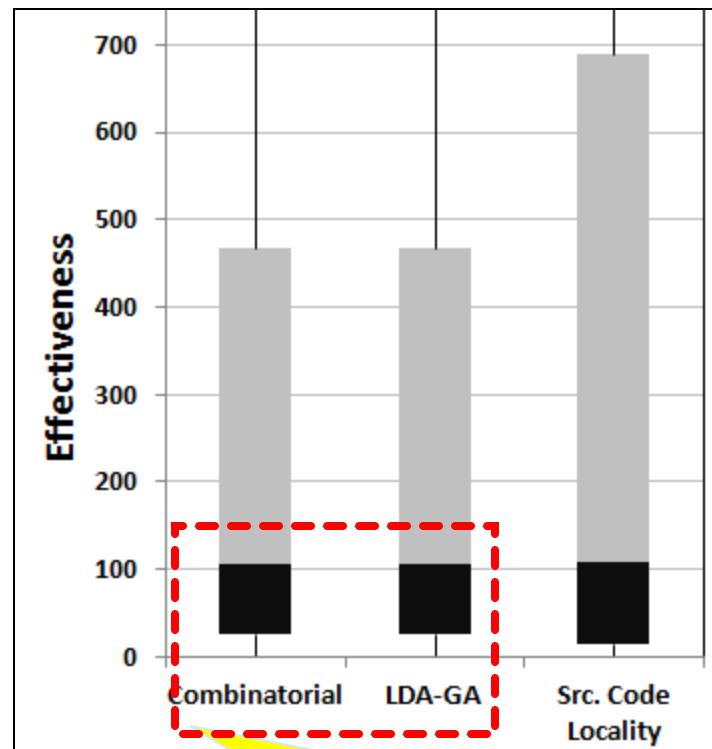


jEdit

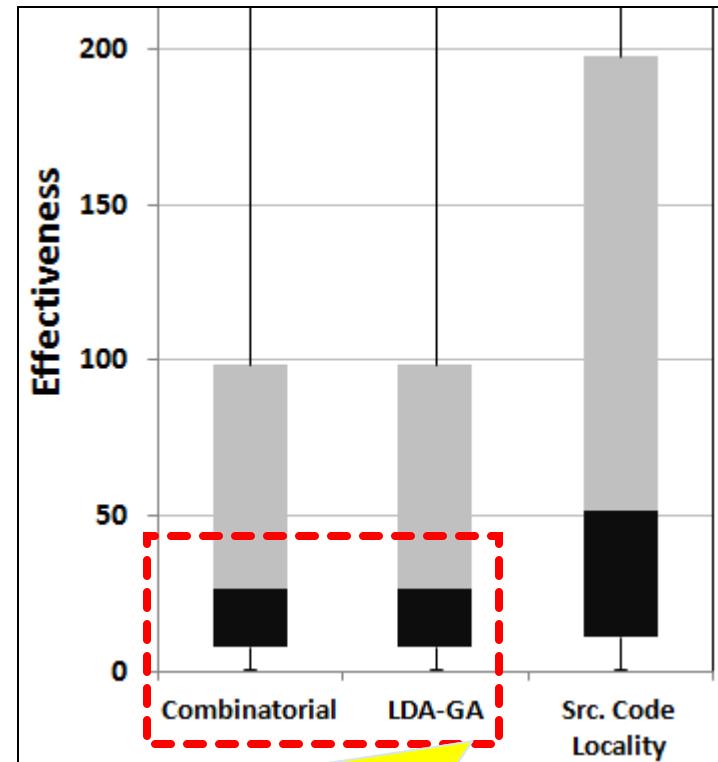


Effectiveness measure

ArgoUML



jEdit

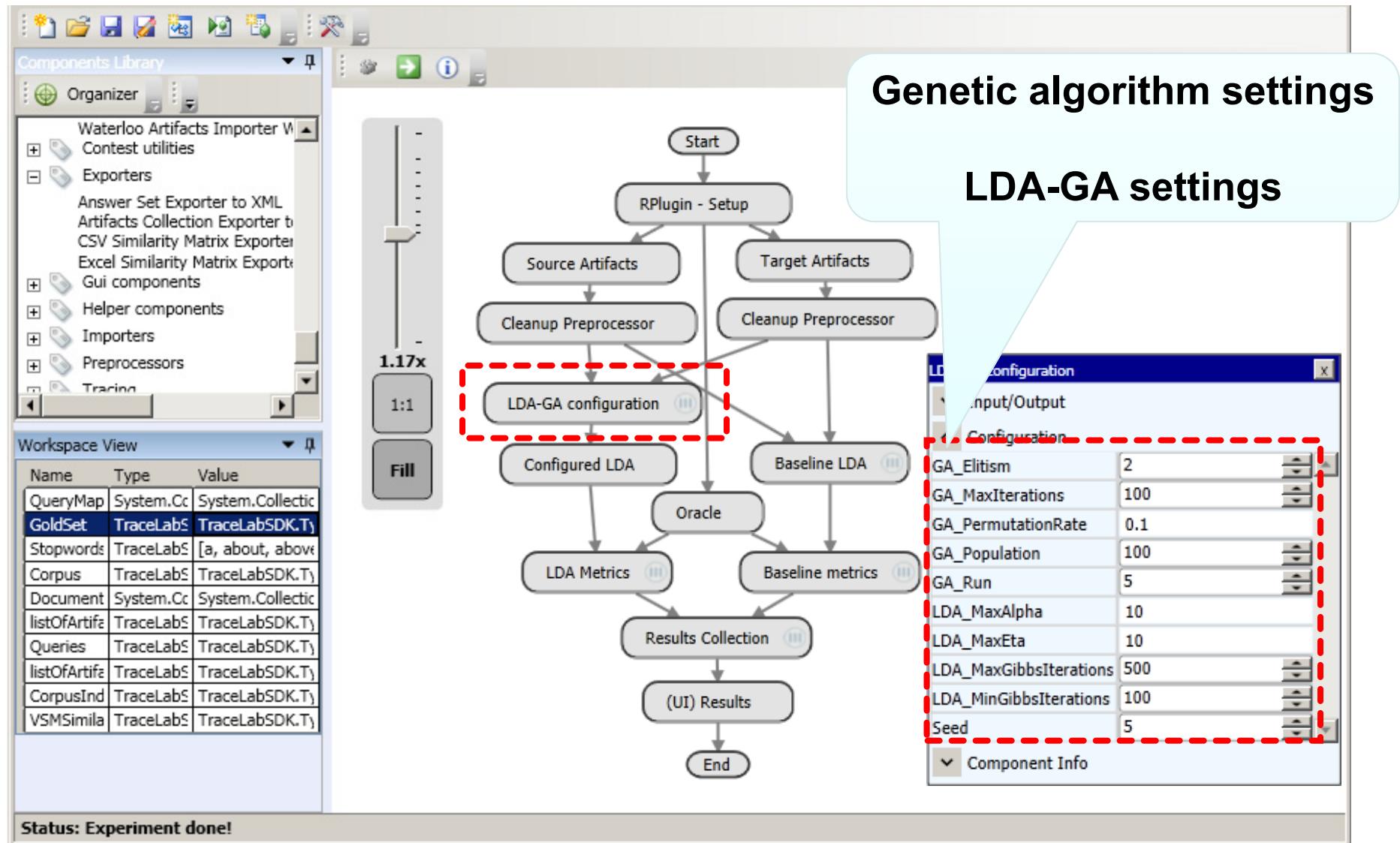


LDA-GA = Combinatorial
LDA-GA outperforms baseline [$p\text{-value} < 0.05$]

Conclusions

- Showed the impact of setting the LDA parameters on the results
- We proposed LDA-GA, a genetic based approach to automatically configure and find the near-optimal solution for LDA parameters
 - Dataset dependent
 - Oracle & task independent
- The approach was evaluated on three maintenance tasks

LDA-GA in TraceLab



Thank you! Questions?

<http://www.distat.unimol.it/reports/LDA-GA/>

<http://www.cs.wm.edu/semeru/data/tefse13/>



UNIVERSITÀ
DEGLI STUDI
DEL MOLISE



References

- [Hindle et al., ICSE'12] A. Hindle, E. T. Barr, Z. Su, M. Gabel, and P. T. Devanbu, "On the naturalness of software," in Proc. of the 34th IEEE/ACM International Conference on Software Engineering (ICSE'12), Zurich, Switzerland, June 2-9, 2012, pp. 837–847.
- [ICPC'10] R. Oliveto, M. Gethers, D. Poshyvanyk, and A. De Lucia, "On the equivalence of information retrieval methods for automated traceability link recovery," in Proc of the 18th IEEE International Conference on Program Comprehension (ICPC'10), Braga, Portugal, 2010, pp. 68–71.
- [SCAM'10] S. Grant and J. R. Cordy, "Estimating the optimal number of latent concepts in source code analysis," in Proc. of the 10th International Working Conference on Source Code Analysis and Manipulation (SCAM'10), 2010, pp. 65–74.

Threats to Validity

- We used datasets that have been used in other studies
- We ran GA 30 times to account for randomness
- Non-parametric statistical test
- Generalizability of results to other SE tasks

GA Settings

- Implementation: GA library in R
- Population size: 100
- Elitism of 2 individuals
- Roulette wheel selection operator
- Crossover probability: 0.6
- Mutation probability: 0.01
- Stop criteria:
 - No improvement in 10 generations
 - When reaching 100 generations

Software Artifact Labeling

TABLE IV
AVERAGE OVERLAP BETWEEN AUTOMATIC AND MANUAL LABELING.

| exVantage | | | | | |
|----------------------|--------|---------------|----------------------|---------|-------|
| | LDA | | De Lucia et al. [13] | | |
| | LDA-GA | Combinatorial | n = M | n = M/2 | n = 2 |
| Max | 100% | 100% | 100% | 100% | 100% |
| 3rd Quartile | 95% | 95% | 71% | 70% | 69% |
| Median | 67% | 70% | 59% | 60% | 54% |
| 2nd Quartile | 60% | 67% | 34% | 50% | 41% |
| Min | 50% | 50% | 0% | 0% | 40% |
| Mean | 74% | 77% | 52% | 56% | 60% |
| St. Deviation | 19% | 17% | 31% | 34% | 23% |

| JHotDraw | | | | | |
|----------------------|--------|---------------|----------------------|---------|-------|
| | LDA | | De Lucia et al. [13] | | |
| | LDA-GA | Combinatorial | n = M | n = M/2 | n = 2 |
| Max | 100% | 100% | 100% | 100% | 100% |
| 3 Quartile | 81% | 82% | 73% | 70% | 66% |
| Median | 71% | 75% | 65% | 61% | 56% |
| 2 Quartile | 47% | 50% | 46% | 45% | 41% |
| Min | 14% | 14% | 0% | 38% | 29% |
| Mean | 65% | 66% | 59% | 60% | 59% |
| St. Deviation | 28% | 26% | 28% | 20% | 24% |