# Topic Modeling

# Topic Models

According to David Blei:

"Topic models are a suite of algorithms that uncover the hidden thematic structure in document collections. These algorithms help us develop new ways to search, browse and summarize large archives of texts"

(http://www.cs.columbia.edu/~blei/topicmodeling.html)

# Topic Models

Topics

Documents

Topic proportions and assignments



(Blei, 2012)

# LDA Topic Models

LDA = Latent Dirichlet Allocation

- ► a topic is a distribution of probabilities of words
- ► all words in a document can belong to all topics
- ► a document is a distribution of probabilities of topics

# LDA Topic Models

*a topic:*                                    *a word:*                                    *a*

| sole (10.1%)<br>cuore (6.4%)<br>amore (4.7%)<br>… |

amore

4.7%     7.1%     5.8%     12.4%     5.2%     15.8%

bad poetry

sentiments

very bad poetry

# LDA: How Does it Work?

► Initialize topic assignments randomly

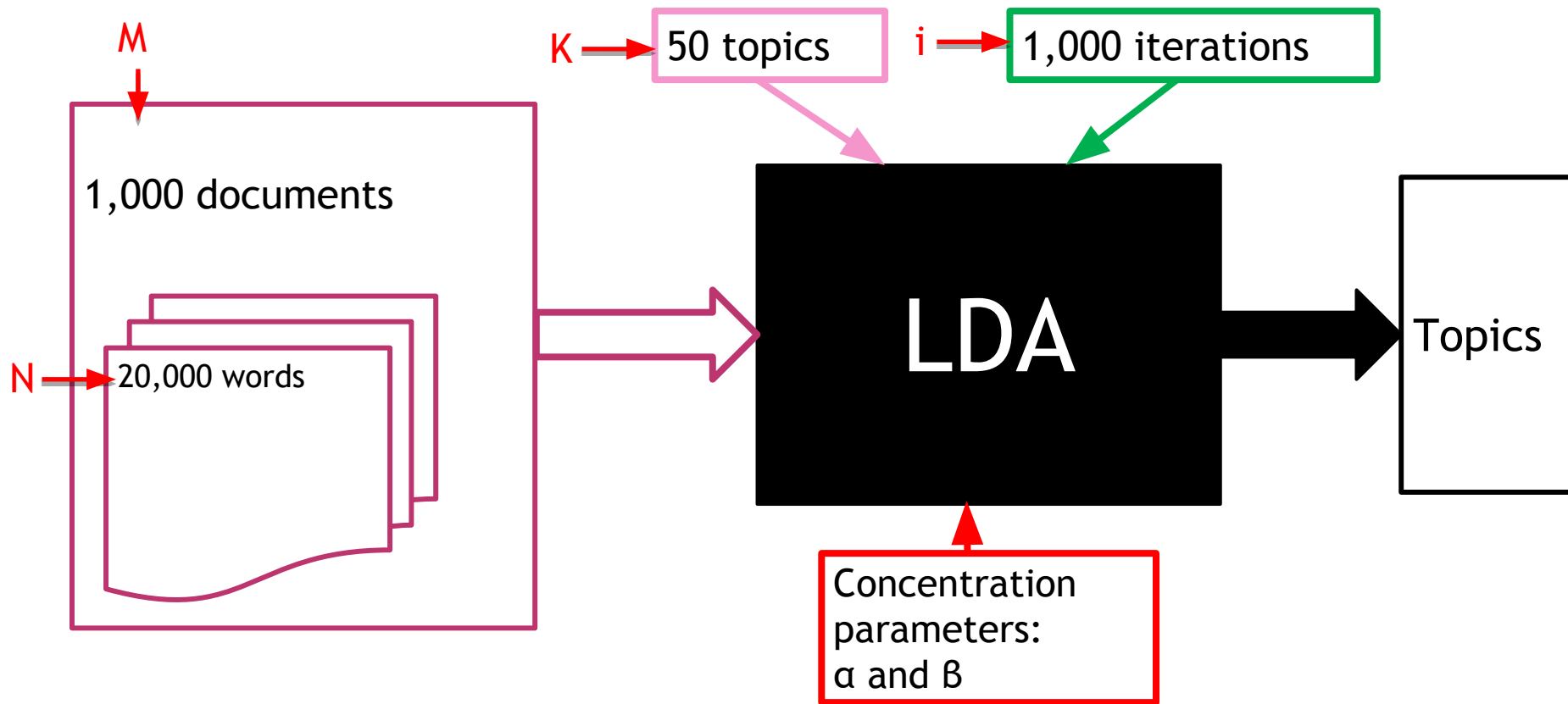► For each word in each document:

  ► re-sample topic for word,
given all other words and their current topic assignments

► Iterate $n$ times!

# LDA: How Does it Work?
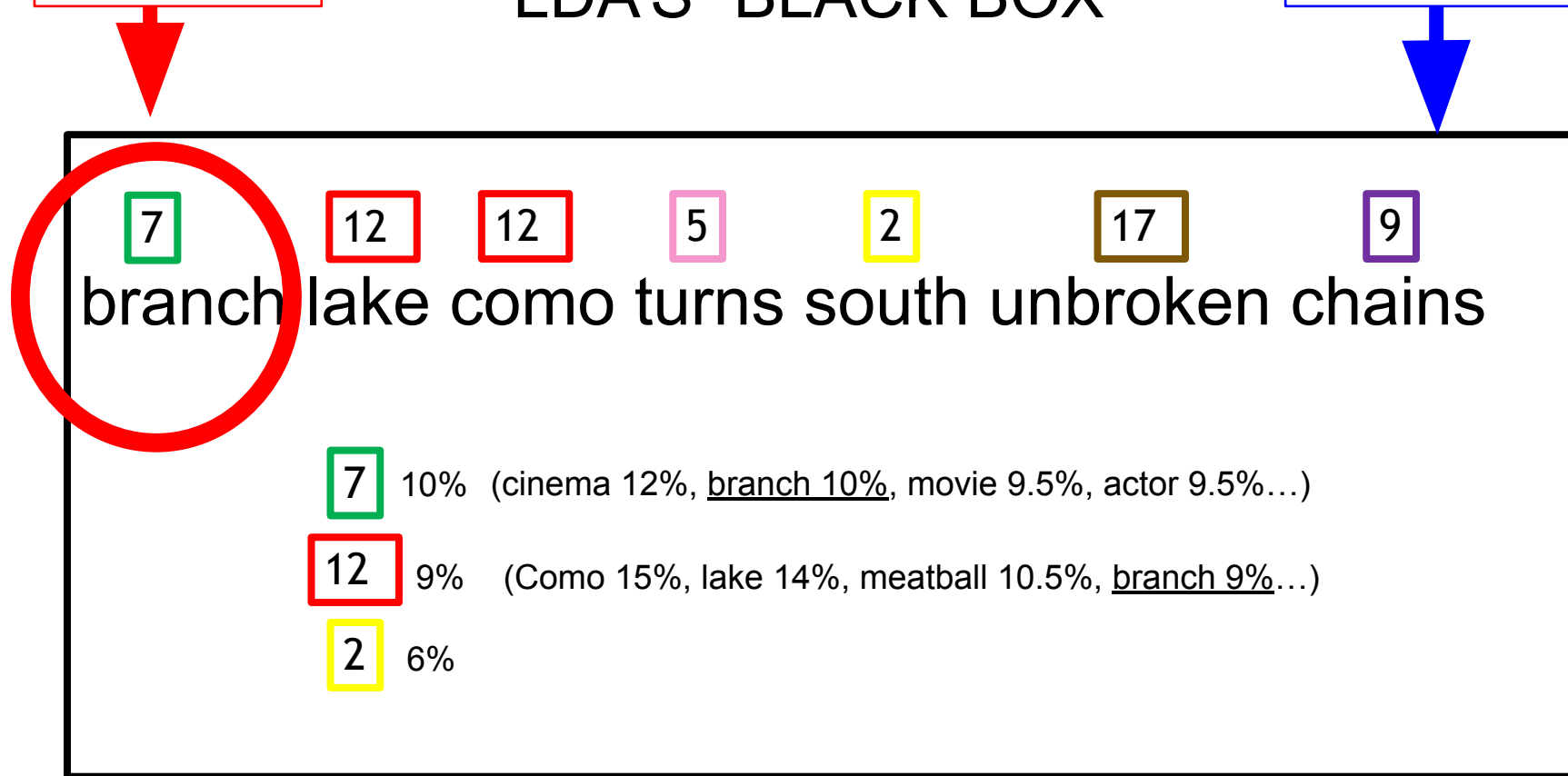
# LDA'S "BLACK BOX"

iteration #1,456

document #151

7 12 12 5 2 17 9

branch lake como turns south unbroken chains

7 10% (cinema 12%, branch 10%, movie 9.5%, actor 9.5%…)

12 9% (Como 15%, lake 14%, meatball 10.5%, branch 9%…)

2 6%

# LDA'S "BLACK BOX"

12   12   5   2   17   9

branch lake como turns south unbroken chains

7   10%   (cinema 12%, branch 10%, movie 9.5%, actor 9.5%…)

12   9%   (Como 15%, lake 14%, meatball 10.5%, branch 9%…)

2   6%