

# Automated Topic Naming to Support Cross-project Analysis of Software Maintenance Activities

Abram Hindle

Dept. of Computer Science  
University of California, Davis  
Davis, CA, USA  
abram@softwareprocess.es

Michael W. Godfrey

David Cheriton School of  
Computer Science  
University of Waterloo  
Waterloo, Ontario, CANADA  
migod@uwaterloo.ca

Neil A. Ernst

Dept. of Computer Science  
University of Toronto  
Toronto, Ontario, CANADA  
nernst@cs.toronto.edu

John Mylopoulos

Dept. Information Eng. and  
Computer Science  
University of Trento  
Trento, ITALY  
jm@disi.unitn.it

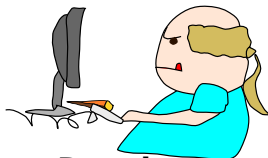
# Who Cares About Quality?



**New Developers**



**Managers**



**Developers**



**Investors**



**Customers**

# What is this commit about?



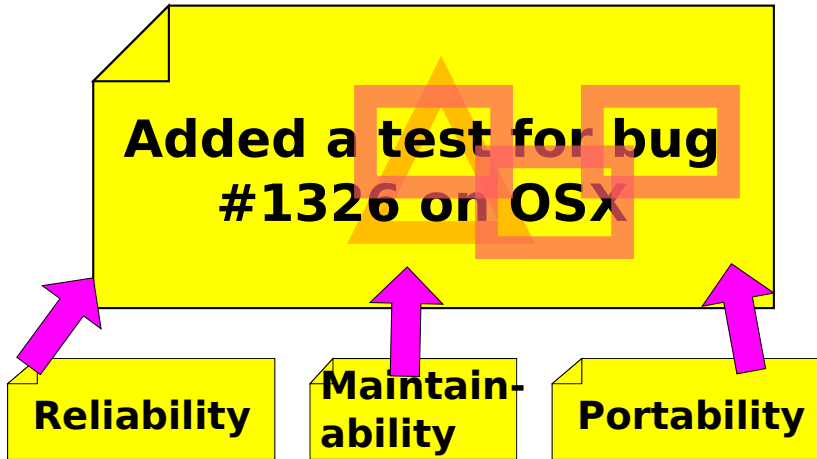
**Added a test for bug  
#1326 on OSX**

# What is this commit about?

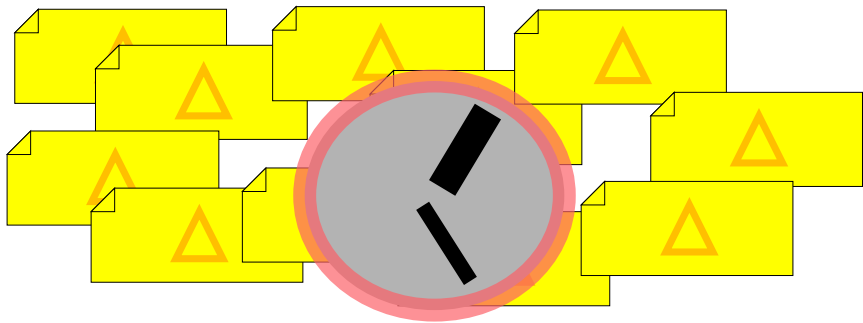


**Added a test for bug  
#1326 on OSX**

# What is this commit about?



# But we have many commits..



**Reliability**

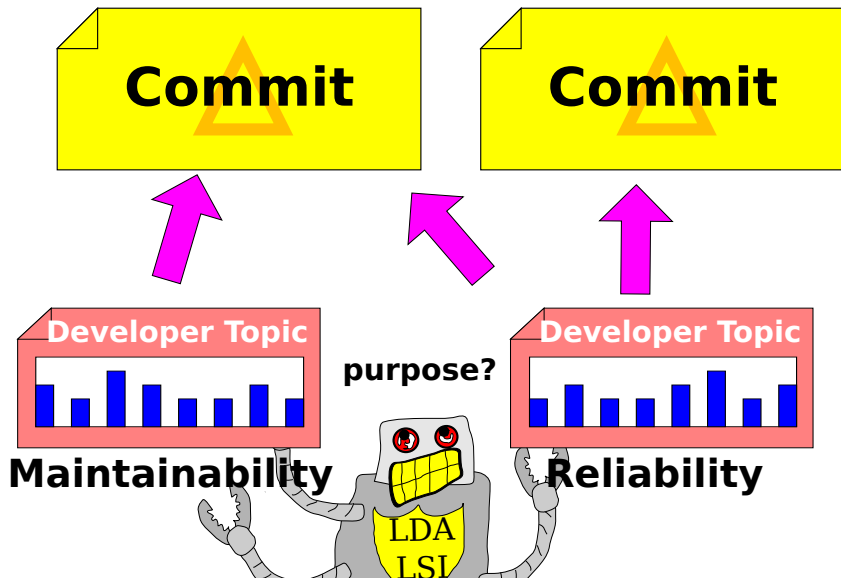


**Maintain-  
ability**

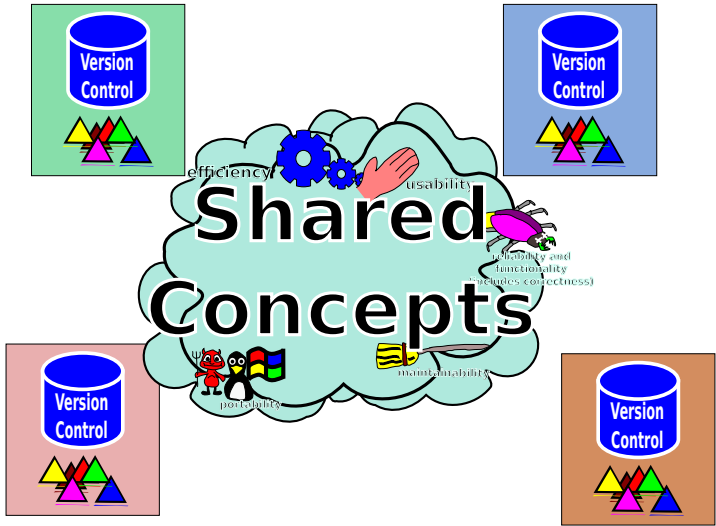


**Portability**

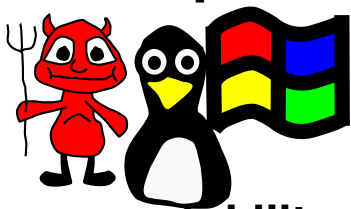
# Developer Topics



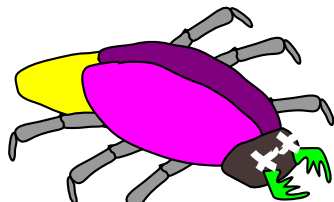
# Cross Project Relevance



# Quality-related Non Functional Requirements (NFRs)



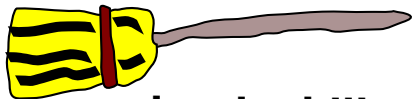
**portability**



**reliability and  
functionality  
(includes correctness)**

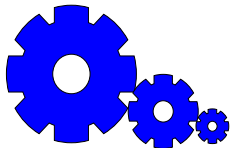


**usability**



**maintainability**

**efficiency**

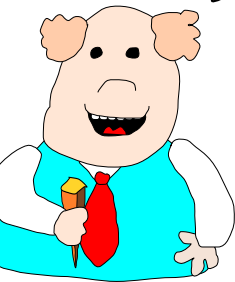


[iso9126]

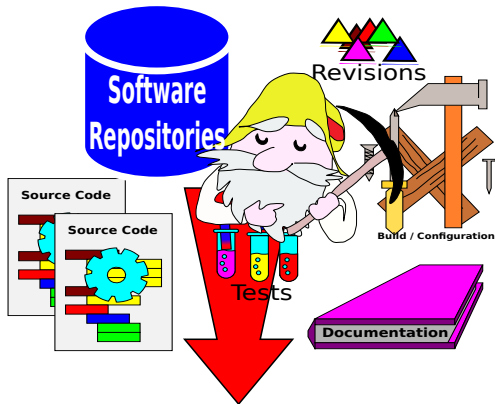
[cleland-huang03]

[ernst19]

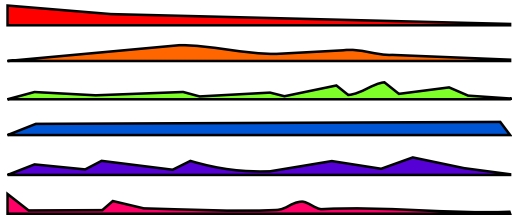
Can't we just summarize quality related efforts within this project?



Maintainability  
Functionality  
Portability  
Efficiency  
Usability  
Reliability



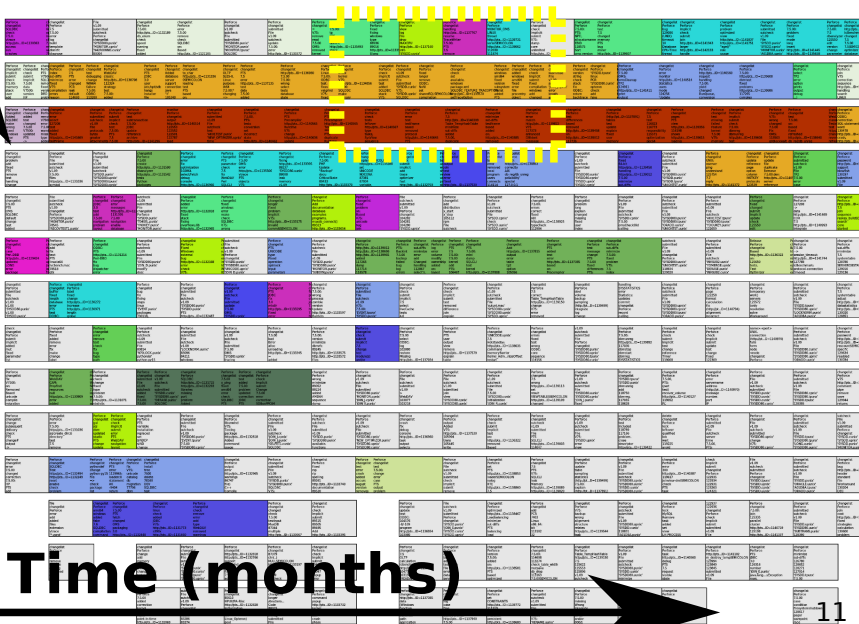
### *Non-Functional Requirements*



time ->

# Labelled Developer Topics

2004 Jun 2004 Jul 2004 Aug 2004 Sep 2004 Oct 2004 Nov 2004 Dec 2005 Jan 2005 Jun 2005 Jul 2005 Aug 2005 Oct 2005 Nov 2005 Dec 2006 Jan 2006 Feb 2006 Mar 2006 Apr 2006 May 2006



Time (months)

# Labelled Developer Topics

Unique Topics

|  |  |  |   |   |
|--|--|--|---|---|
| changelist<br>Perforce<br>7.5.00:<br>tp<br>PTS<br>size<br>dbmrfc<br>http://pts...ID=1135493<br>stack<br>Perl | changelist<br>Perforce<br>fixed<br>fixing<br>windows<br>type<br>88936<br>89016<br>http://p...ID=1135493<br>Ullan | changelist<br>Perforce<br>view<br>log<br>file<br>drop<br>MAXCPU<br>http://pts...ID=1137160 | Perforce<br>changelist<br>handling<br>http://pts...ID=1137767<br>resume<br>TraceWriter<br>7.5.00:<br>PTS<br>7.5<br>remove | changelist<br>Perforce<br>MONITOR_OMS<br>LINUX<br>thread<br>http://pts...ID=1138731<br>serialSEMICOLON<br>http://pts...ID=1138662<br>112030<br>111874 |
|--|--|--|---|---|

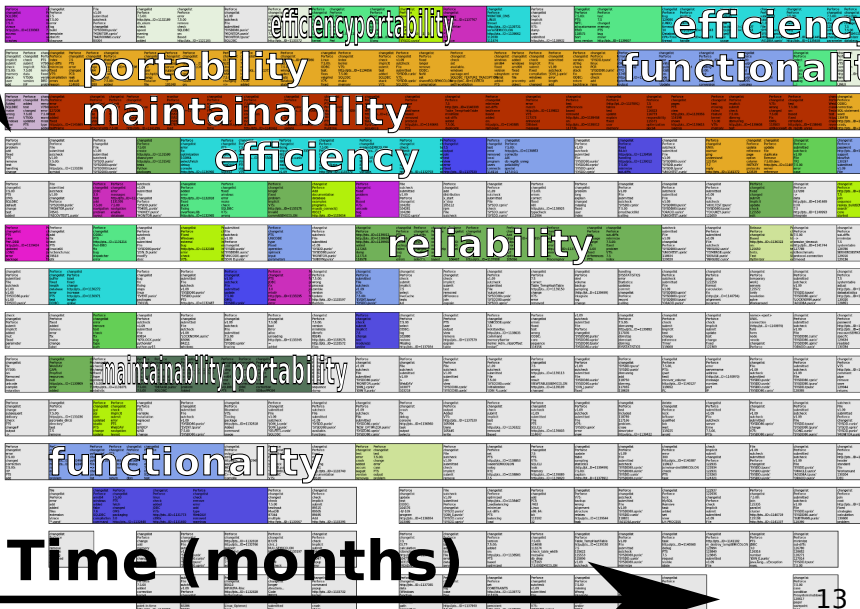
|   |  |   |  |   |
|---|--|---|--|---|
| changelist<br>Perforce<br>Linux<br>7.5.00:<br>ODBC<br>test:<br>fixes<br>valgrind<br>V75:<br>fixed | Perforce<br>changelist<br>kernel<br>http://pts...ID=1134856<br>7.5.00:<br>SQLDBC<br>make<br>memory | changelist<br>Perforce<br>submit<br>test<br>test<br>'SYSDBA unix'<br>st k<br>V... | Perforce<br>changelist<br>check<br>set<br>available<br>window<br>SQLDBC FEATURE TRACOPTION<br>LOMhttp://pts...ID=1133570<br>setTraceOption | Perforce<br>changelist<br>windows<br>amd64<br>version<br>subsys<br>added<br>PTS<br>error<br>file<br>procentry<br>object |
|---|--|---|--|---|

|   |  |   |   |   |
|---|--|---|---|---|
| changelist<br>Perforce<br>PTS<br>Precompiler<br>http://pts...ID=1140565<br>Tests<br>***<br>PTS:<br>duplicate<br>Runtime | changelist<br>Perforce<br>check<br>tests<br>leftover<br>http://pts...ID=1140567<br>type<br>tasks,<br>SysView<br>sequence | changelist<br>Perforce<br>test<br>testframe<br>based<br>removed<br>protocol<br>http://pts...ID=1140645<br>PTS | changelist<br>Perforce<br>7.5<br>http://pts...ID=1140309<br>Table_TempHashTable<br>minimize<br>sut-diffs<br>added<br>7.5.00,<br>duplicate | Perforce<br>changelist<br>minimize<br>sut-diffs<br>testframe<br>based<br>added<br>on:<br>http://pts...ID=1140385<br>removed |
|---|--|---|---|---|

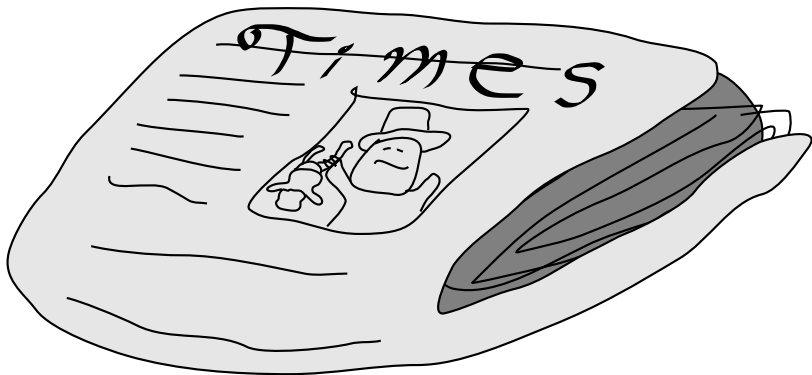
Time (months)

# Labelled Developer Topics

2004 Jun 2004 Jul 2004 Aug 2004 Sep 2004 Oct 2004 Nov 2004 Dec 2005 Jan 2005 Jun 2005 Jul 2005 Aug 2005 Oct 2005 Nov 2005 Dec 2006 Jan 2006 Feb 2006 Mar 2006 Apr 2006 May 2006

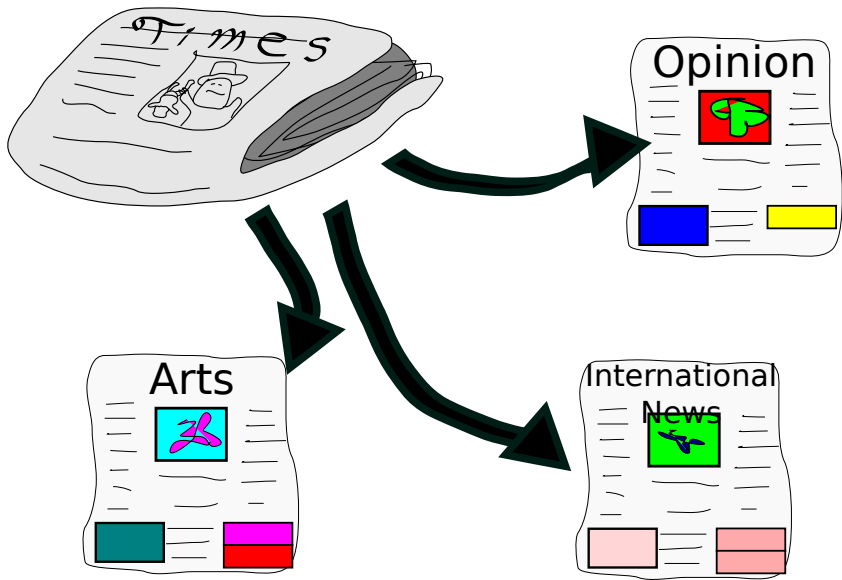


# Example



[Blei]

**apologies to those with  
prior LDA/LSI experience**

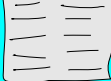


## Arts

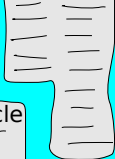
### Section



Article



Article



Article



## International

### News

### Section



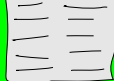
Article



Article



Article



Article



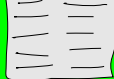
Article



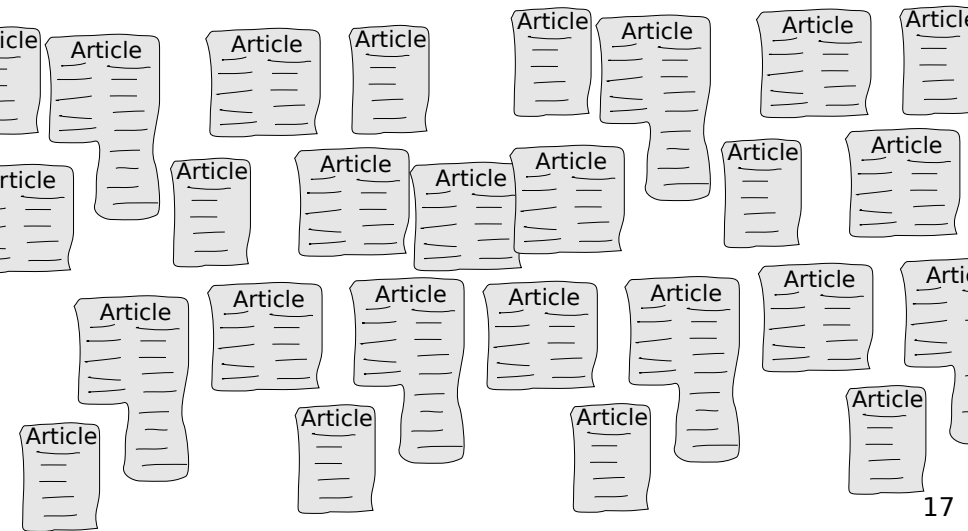
Article

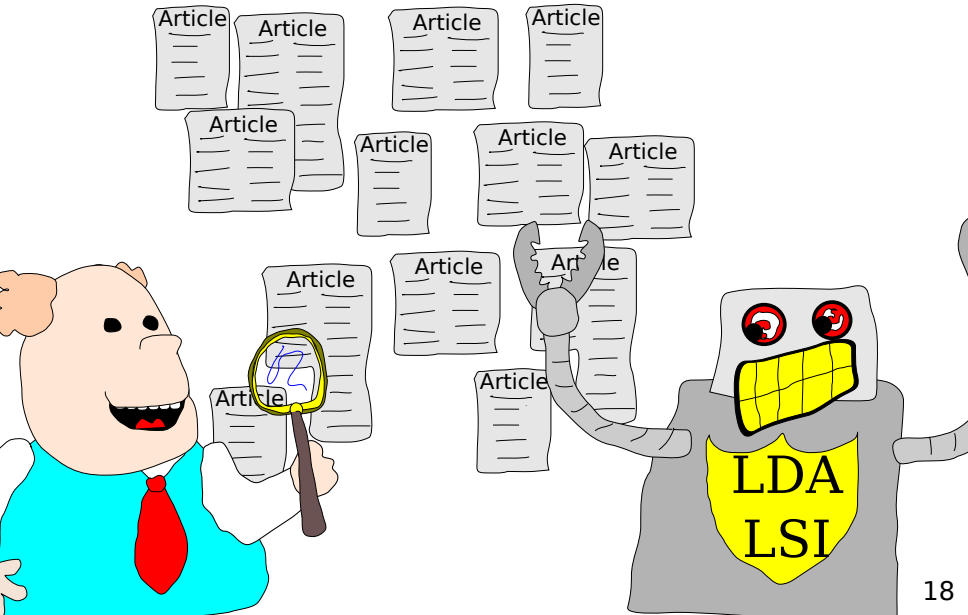


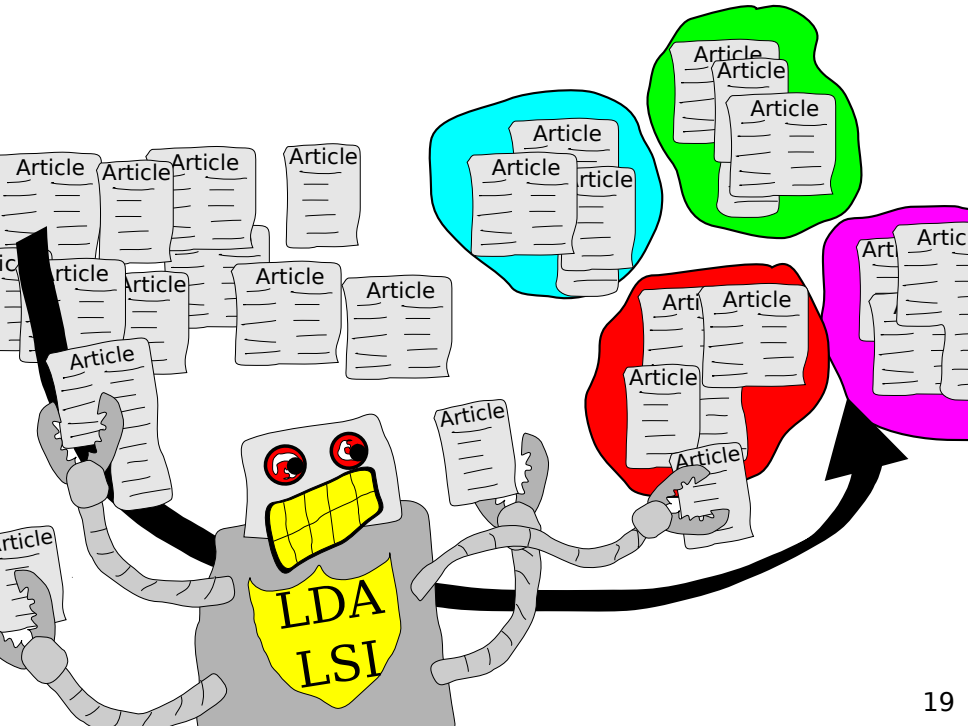
Article

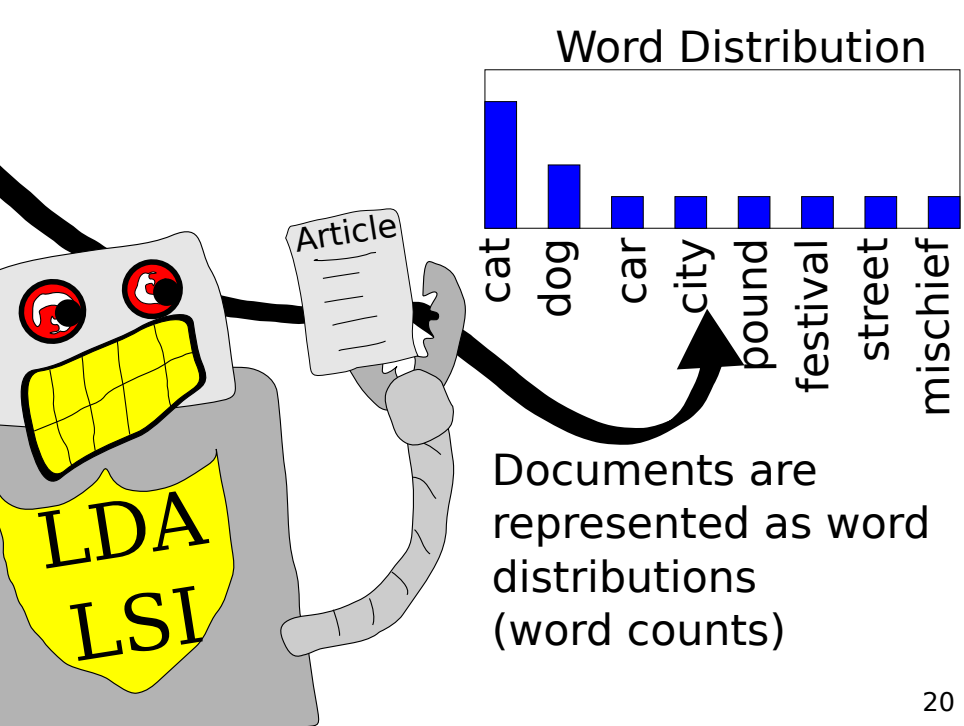


What if we didn't know what section the articles were in?

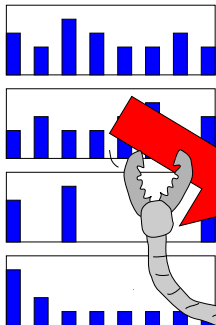




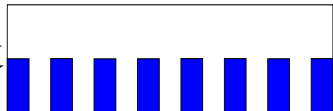




## Word Distributions

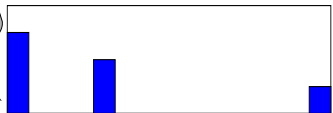


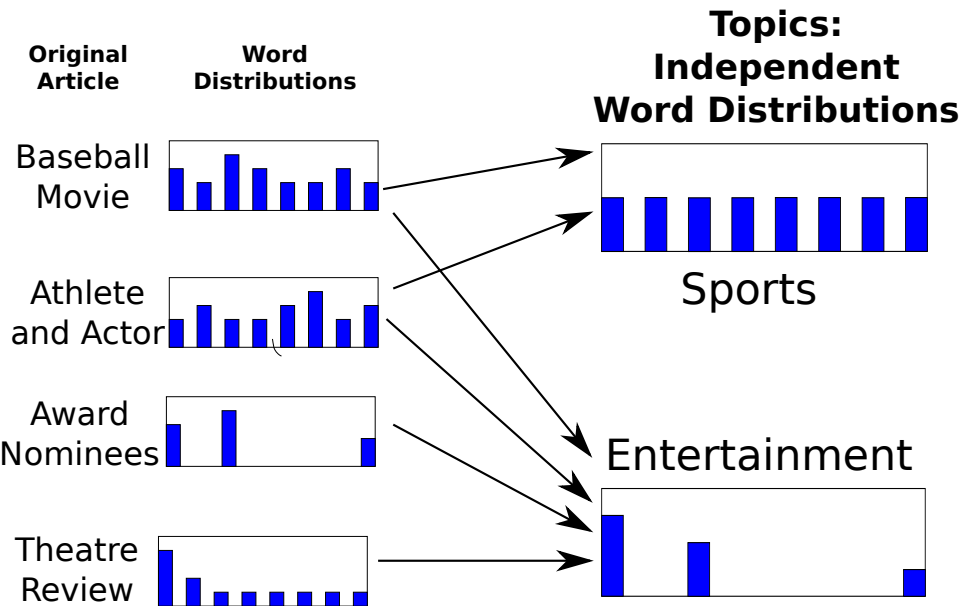
## Topics: Independent Word Distributions



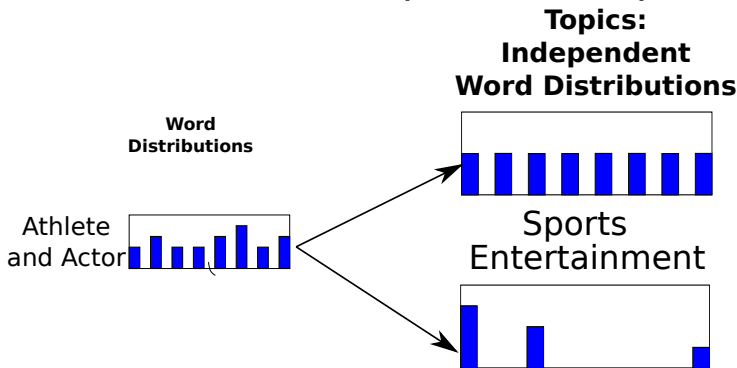
LDA finds independent word distributions that the documents are related to.

Documents can be associated with more than one topic.



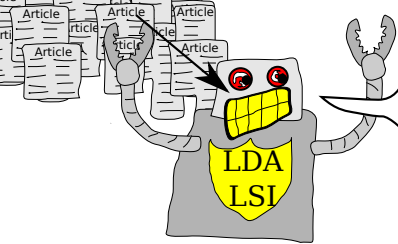


# Documents are represented as a linear combination of independent topics



$$\begin{array}{l} C_0 \times \text{[Sports Distribution]} \\ + \\ C_1 \times \text{[Entertainment Distribution]} \end{array} \approx \text{[Document Distribution]}$$

This equation shows the linear combination of two topics to form a document. The first term,  $C_0$  multiplied by the 'Sports' topic distribution, and the second term,  $C_1$  multiplied by the 'Entertainment' topic distribution, are added together. The result is approximately equal to the document's word distribution, which is shown as a bar chart on the right. The document distribution is a weighted sum of the two topic distributions, where the weights are  $C_0$  and  $C_1$ .



Here are two topics. I don't know what they are about!

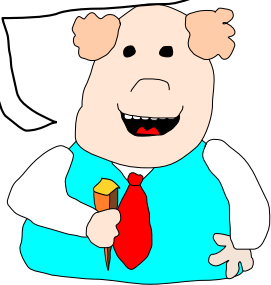
## Topic 1

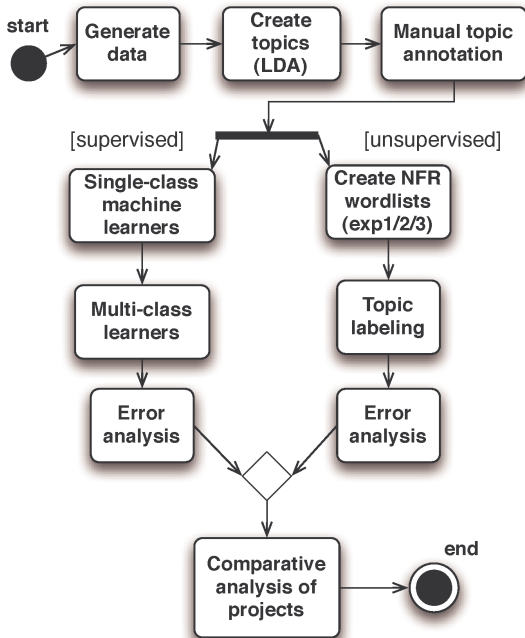
- \* play
- \* game
- \* inning
- \* player
- \* quarter
- \* opponent
- \* ...

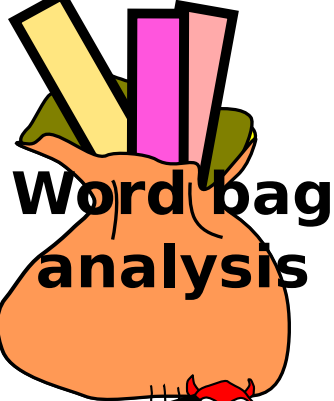
## Topic 2

- \* gambling
- \* play
- \* night life
- \* comedy
- \* movie
- \* theatre
- \* ...

These word lists look look like: **Sports** and **Entertainment** !







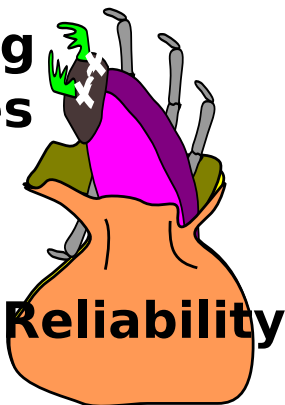
# Word Bag Examples



## Portability

portability  
transferability  
interoperability  
documentation  
internationalization  
i18n

...



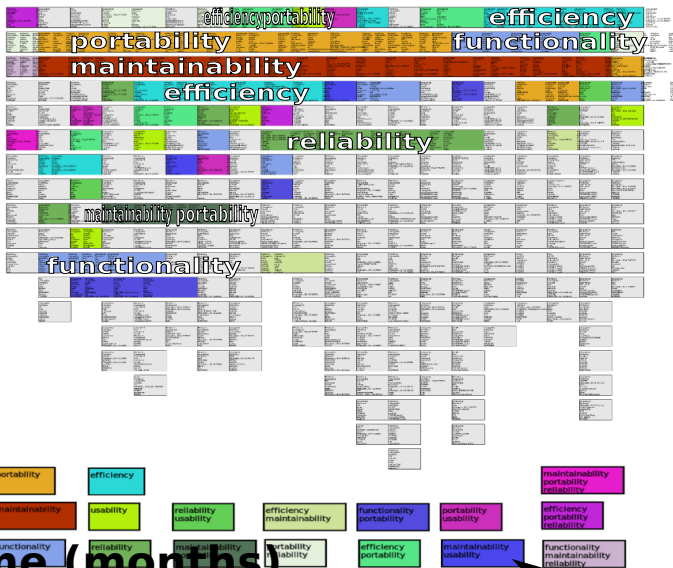
## Reliability

reliability  
failure  
error  
redundancy  
fails  
bug

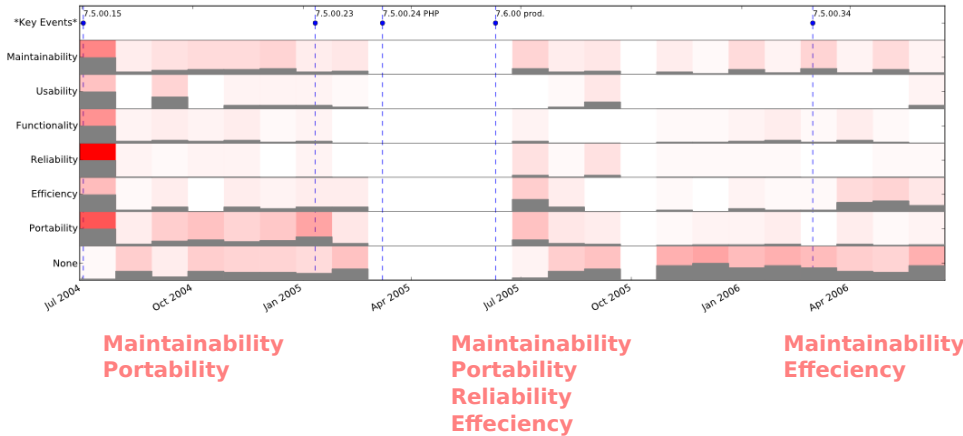
...

# Labelled Topics of MaxDB 7.500

2004 Jun 2004 Jul 2004 Aug 2004 Sep 2004 Oct 2004 Nov 2004 Dec 2004 Jan 2005 Jun 2005 Jul 2005 Aug 2005 Oct 2005 Nov 2005 Dec 2006 Jan 2006 Feb 2006 Mar 2006 Apr 2006 May 2006 Jun 2006



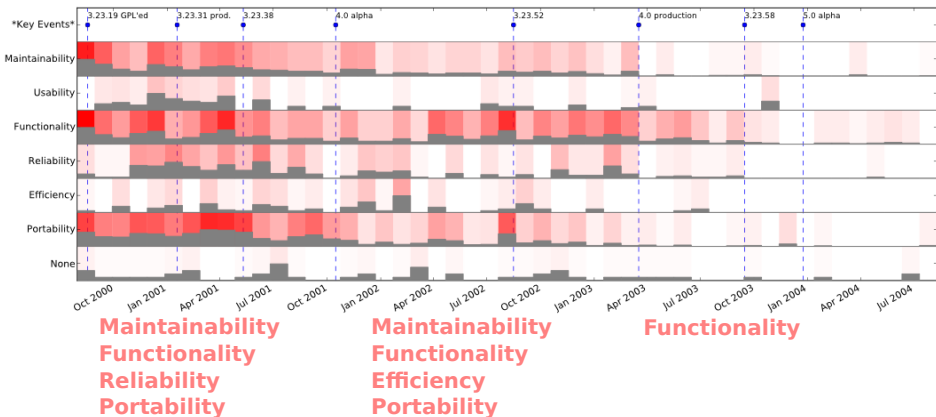
# MaxDB 7.500 Timeline



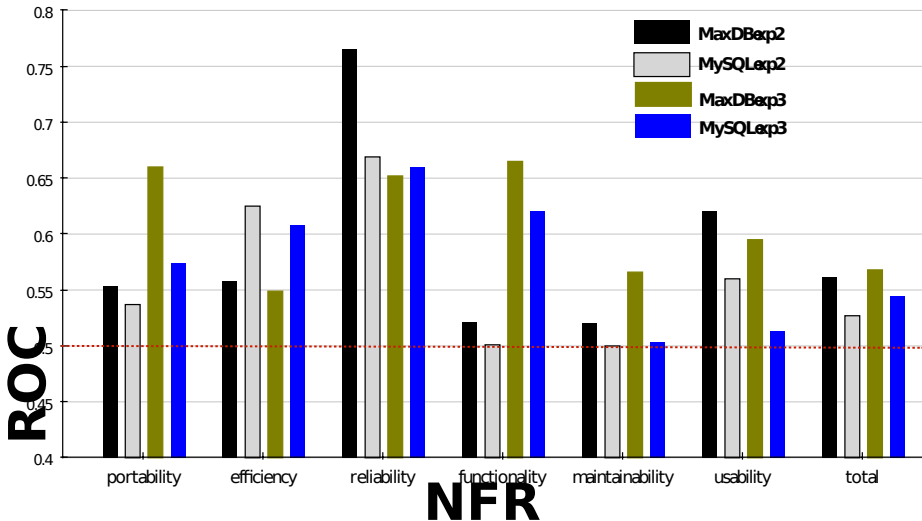
## Unique Topics



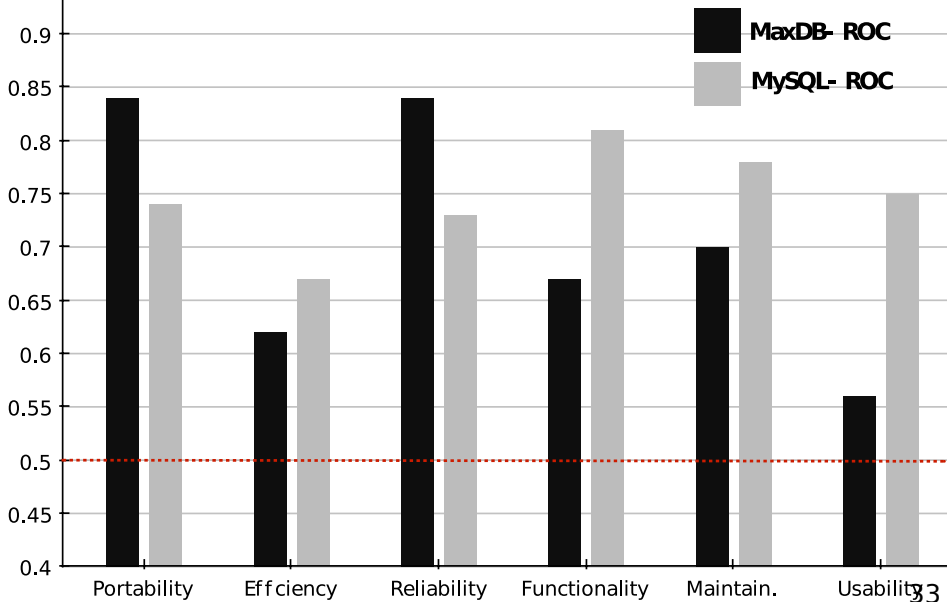
# MySQL 3.23 Timeline



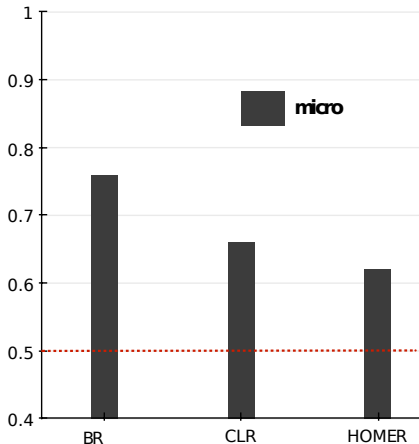
# ROC Values of Semi-Supervised Word Bags



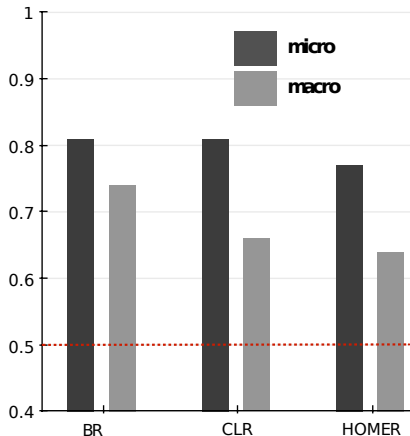
# Supervised Tags



# Supervised Multitag Classifiers: MySQL and MaxDB

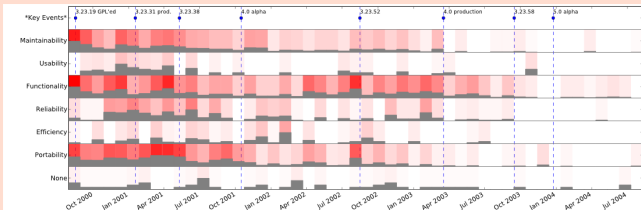
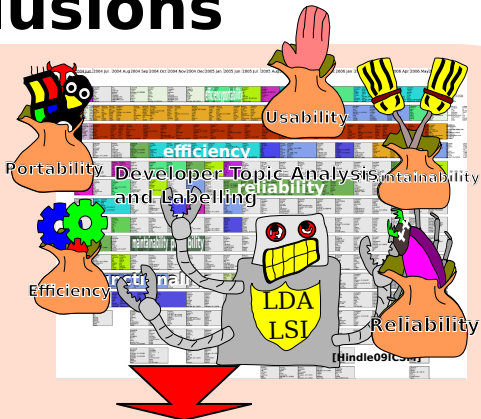
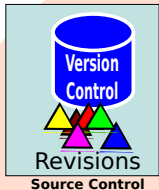


**MySQL  
Classifiers**

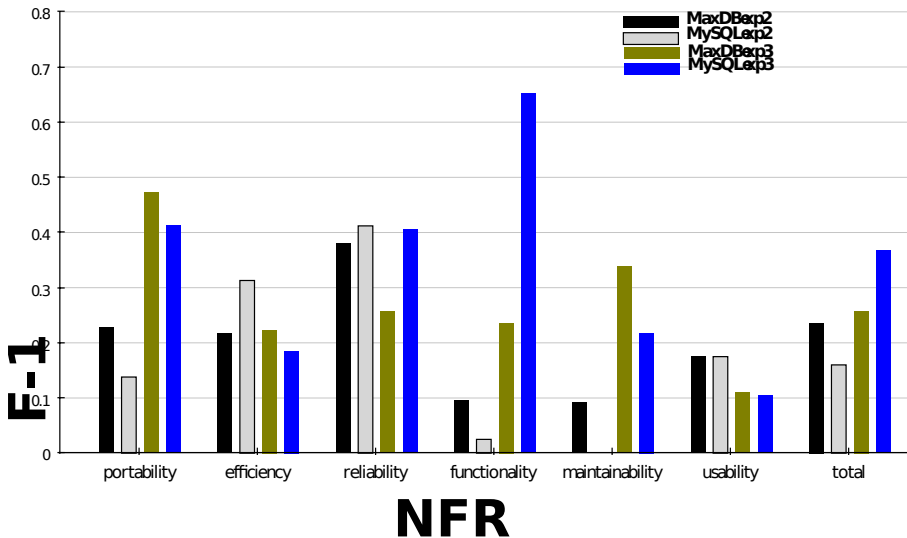


**MaxDB  
Classifiers**

# Conclusions



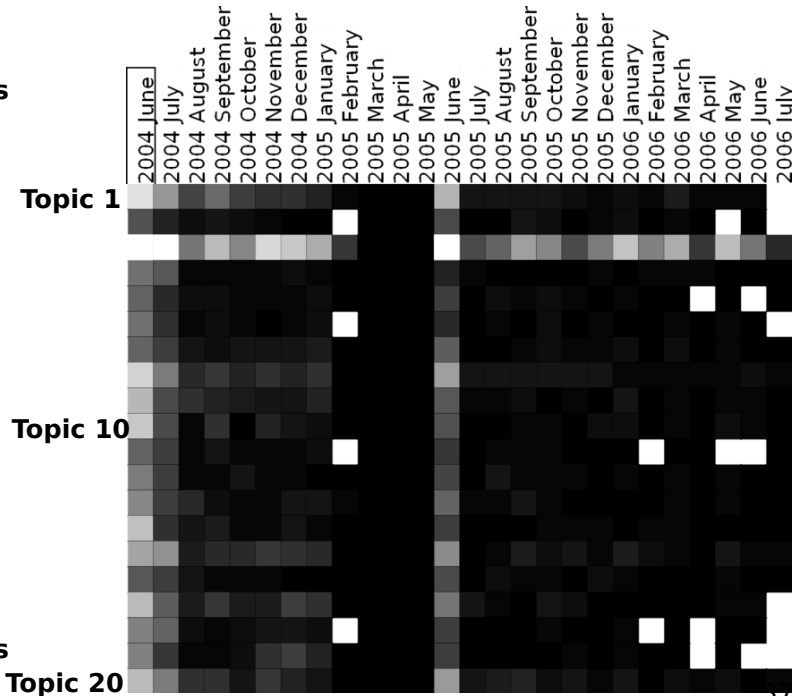
# F-1 Measure of Semi-Supervised Word Bags



**Many Documents**

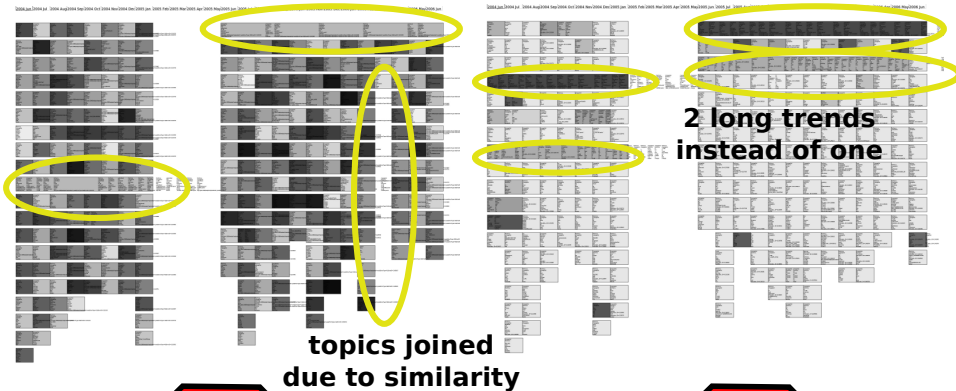


**Few Documents**



# Annotation: Stop Words

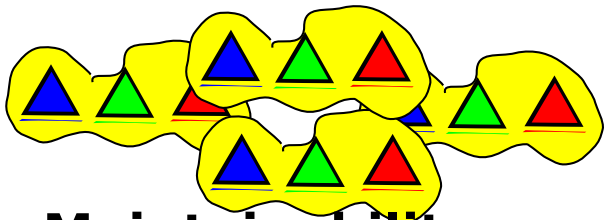
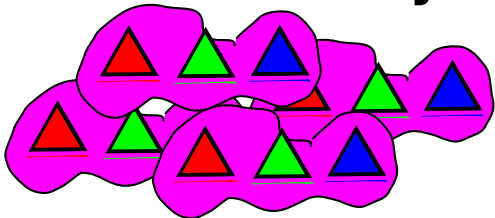
## MaxDB 7.500 Case Study



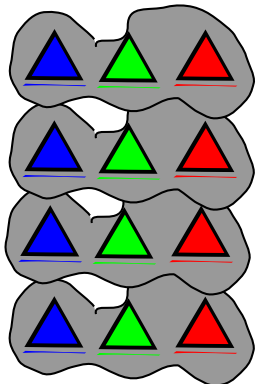
# Annotation: Training Sets



**Maintainability+**



**Maintainability-**



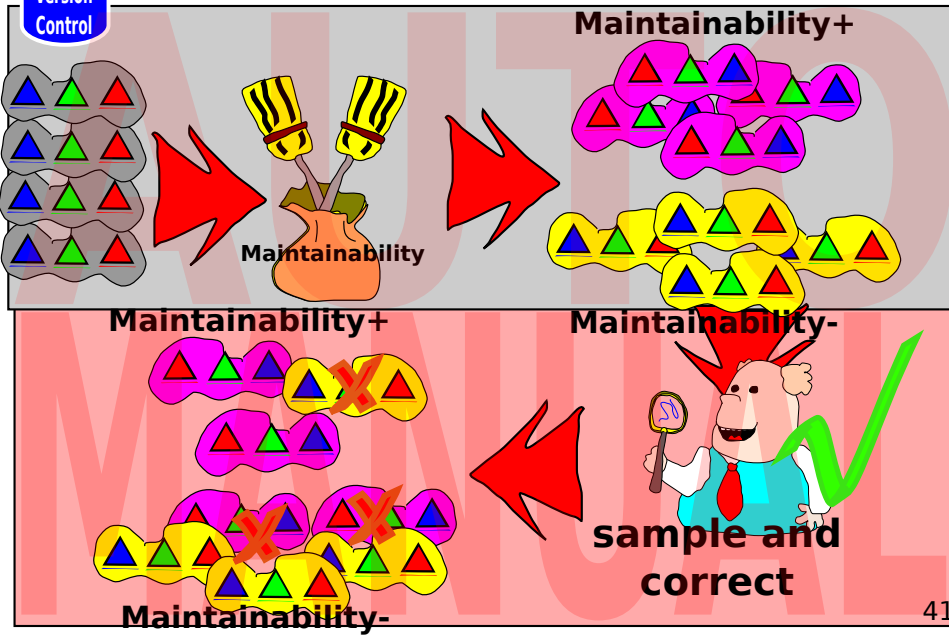
# Annotation: Stop Words



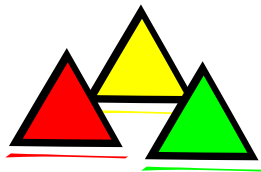
**Used in topic analysis  
or to reduce # of  
features for learners.**

perhaps clearly between  
done there who because  
haven't move in asking  
nevertheless example  
sensible our some  
elsewhere upon ask  
beforehand ie found  
anywhere it containi  
everywhere det  
need associati  
specifying  
con di  
for

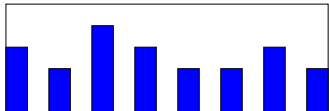
# Annotation: Training Sets



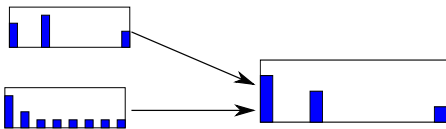
# Message



## Word Distribution



## Topic



**Top 10 Words:**

- \* perforce
- \* bug #
- \* POSIX
- \* Opteron
- \* ...

## Trend

