

# Dealing with large volumes in statistical analysis and text mining tools

April 16-17 2012

II-SDV Nice

*Laurent Hill*  
*Questel*

Confidential



# Introduction

“What’s hot in the semiconductor industry?”

“What’s patented in Korea but not in the US?”

“What are the trends in biofuel production?”



To answer those questions you need to analyze 10,000s to millions of patent records in detail

# Presentation Outline

- Problem statement
- Large scale text mining
- Text mining applications in Fampat
- Statistical analysis at work on large data –sets
- Converging scale-out and small corpus techniques

# Scaling out patent analysis: Problem statement

- Some projects deal with up to ~ 10,000 patents
  - Existing tools well adapted
  - Classic cycle: extract, run text mining, vet data, graph results
- For some other problems, the volume is just too high for this approach
  - See introduction questions
    - Last 10 years of semiconductor patents: 854,000 families, 1M+ patents
    - Patented in Korea but not in the US, 2000-2010: 1,134,000 families
    - Trends in biofuel production: 18,800 patent families (51,000 patents)
  - Yet we need fast answers, stakeholders do not care about volume, they just need answers.
  - An entirely different approach is needed, while still having the same user interface, and without giving up on any functionality for smaller data-sets.

# Large scale text mining

- Text mining:
  - Producing a model from text analysis
  - Deriving high quality information / features from the model
- Issues
  - Text mining analysis is resource intensive, cannot be done in real time on a large corpus
  - The model often is a black box to the user, lowering user confidence on derived applications

# Concept extraction

- Main idea: extract concepts rather than words as a semantic model of patent documents, weight concepts according to sentence type
  - Weighted concepts produce a natural semantic vector model of the patent
- To deal with large volume we need to do this upfront when loading the data
  - 3 CPU years to pre-process the whole database

# Key sentences tagging

- Identify key sentences describing patent object, advantages and drawbacks, independent claims
- Uses morpho-syntactic analysis to spot important sentences
  - “patent writer sentiment analysis”
- Good compromise between conciseness of bibliographic abstracts and full text

# Semantic concept tagging

- « noun phrases » identification
  - Part-of-speech + stemming
  - Verb Suppression, some adjectives
  - Suppress patent boilerplate terminology
    - « preferred embodiment », « skilled artisan » ....
  - Syntactic normalization:
    - « surface of screens » → « screen surface »
- Relevance score computation, based on :
  - Field, key sentence morpho-syntactic detection, and number of occurrences



# Concept extraction example

The invention relates generally to **molecular level cleaning of parts** by **vapor degreasing**.

More particularly, the invention relates to a **solvent mixture** comprising **n-propyl bromide**, a [mixture of **low boiling solvents**] and, ...

The **solvent mixture** of the invention is **non-flammable, non-corrosive, non-hazardous**, and has a **low ozone depletion potential**.

# Normalization examples

the heat conductivity  
the conduction of heat  
heat conductivity  
heat conduction

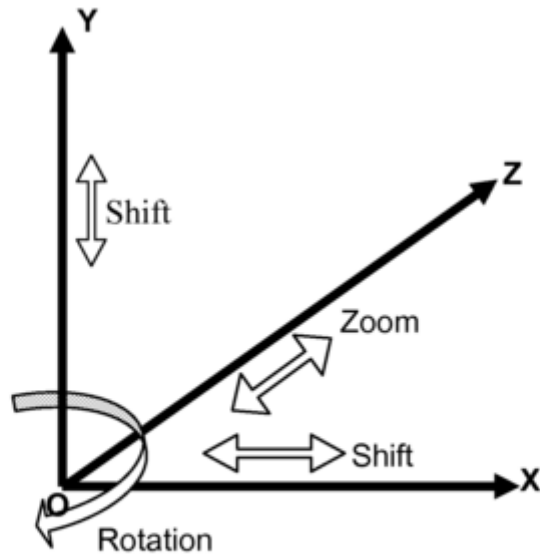


**HEAT CONDUCTION**

the user of the cellular telephone  
cellular phone users  
a user of a mobile telephone  
any mobile phone users  
their user mobile telephones



**MOBILE PHONE USER**



« tag cloud »  
style  
visualization in  
orbit

*Vector model made  
user-friendly*

ALPHANUMERICAL PUSH BUTTON | APPROACH | BLUETOOTHTRADE MARK | BORDERLINE SITUATION |  
**CAMERA GESTURE** | CAMERA MOVEMENT  
**HISTORY** | CAMERA ROTATION **GESTURE** | CAPACITIVE SENSITIVE PAD | CELL PHONE CAMERA |  
 CELLULAR TELEPHONE CAMERA | **CLASSIFIER** | CLASSIFIER ROUTINE | COMMON  
 MOVEMENT DISTANCE | COMPARING POSITION | COMPUTER READABLE INSTRUCTION |  
 COMPUTER READABLE MEDIUM STORING COMPUTER EXECUTABLE  
 INSTRUCTION | CURRENT **GESTURE** | **CURRENT VIDEO IMAGE** | DESKTOP  
 COMPUTER KEYBOARD | DETECTED **GESTURE** | DIGITAL CAMERA | EARLY IMAGE | FEATURE POINT |  
 FEATURE POINT NUMBER | HANDHELD COMPUTER MOUSE | **HISTORICAL GESTURE**  
**INFORMATION** | HISTORY DURATION | HYPERTEXT MARKUP LANGUAGE | IDENTIFIED  
 FEATURE NUMBER | IDENTIFIED FEATURE POINT | IMAGE SUCCESSIVE FRAME | INTEL  
 CORP | LAST **GESTURE** INDICATION | LOCAL DESKTOP COMPUTER | LUCAS KANADE PYRAMIDAL  
 IMPLEMENTATION | MICROPROCESSOR RESEARCH LAB | MINIMUM NUMBER | MOBILE CELLULAR  
 TELEPHONE | MOBILE TELEPHONE KEYPAD | MOVEMENT **GESTURE** LACK | MOVEMENT SPEED |  
 MOVING SPEED | MULTIPURPOSE MICROPROCESSOR | MUSIC PLAYER | NAVIGATION | ONSCREEN  
 CURSOR INTUITIVE MOVEMENT | OUTPUT **GESTURE** | PATTERN RECOGNITION | PRIOR VIDEO IMAGE |  
 RECORDED **GESTURE** HISTORY FILE | RELATIVE POSITIONING | REMOVEABLE FLASH MEMORY | REMOVEABLE  
 MEDIA | **ROTATION GESTURE** | ROTATION **GESTURE** OUTPUTTING | ROTATION

# Text mining applications at database level

- Since text mining has been applied massively, it can be leveraged at the database level (Fampat)
- Similarity searching
  - More like this (one or more patents)
  - Refine by example
- Related concept search
  - “chaussure de ski” yields:
    - Ski boot
    - Skier foot
    - Ski binding

# Large scale statistical analysis

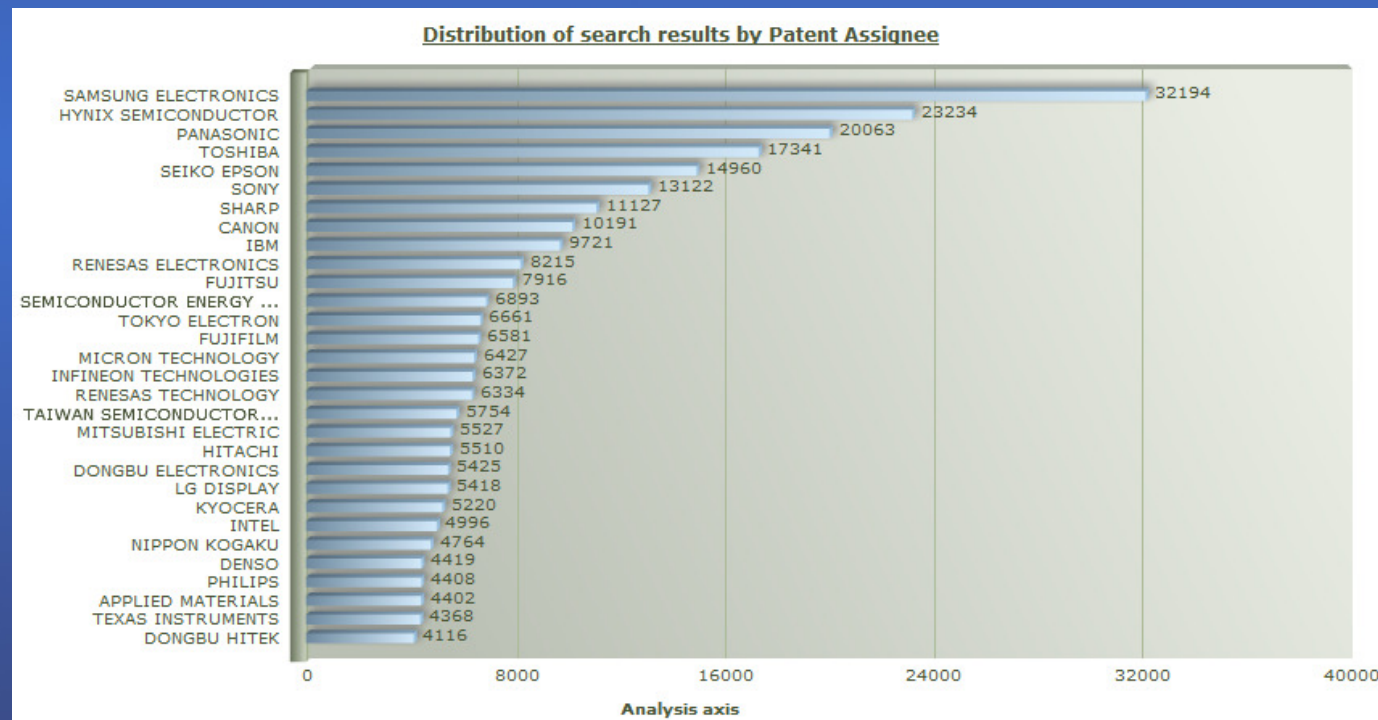
- Nothing special here ?
  - More precisely: that's the goal
  - Be able to get all usual business graphics on patents
  - With close to instant response time on 1 million patent families
  - And *exact results*

# Prerequisites

- All data normalized and clean in the database
  - Normalized assignee names
  - Legal status for all patents and European national phases
  - Normalized text mining concepts
- Everything available at the family level in the same database
- A fast hybrid engine
  - Boolean full-text
  - Semantic
  - Analytic (OLAP)

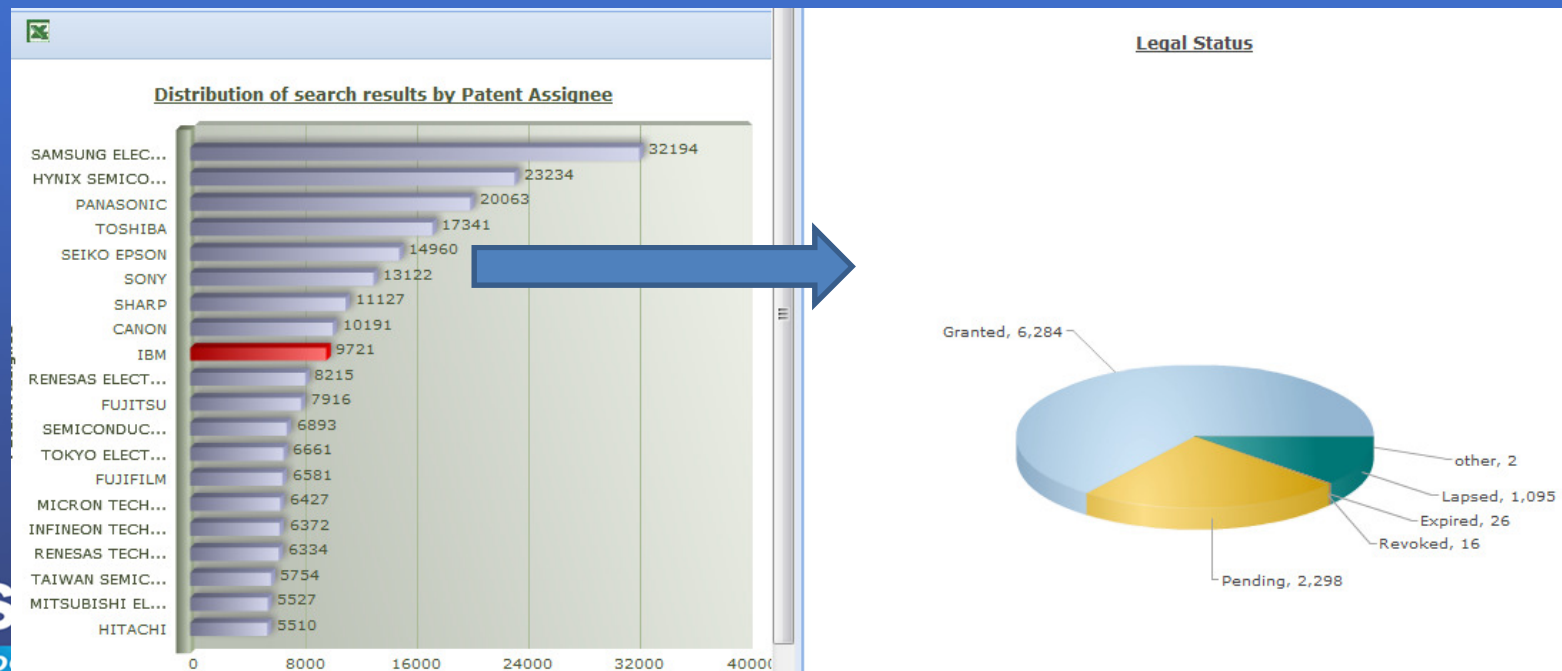
# Results

- 1M patent family, less than 10s for top assignees



# Data exploration

- On large volumes, you often are formulating and testing hypothesis on the spot
  - One click drill down :







# Graphing along Custom axis

- Sometimes you want to produce business graphics on your own queries
  - Combining International classes with concepts
    - A61Q or cosmetic/KEYW
  - Defining precise patent categories by full text queries
    - (bio\_diesel P Alga??)/CLMS

# Custom business graphic

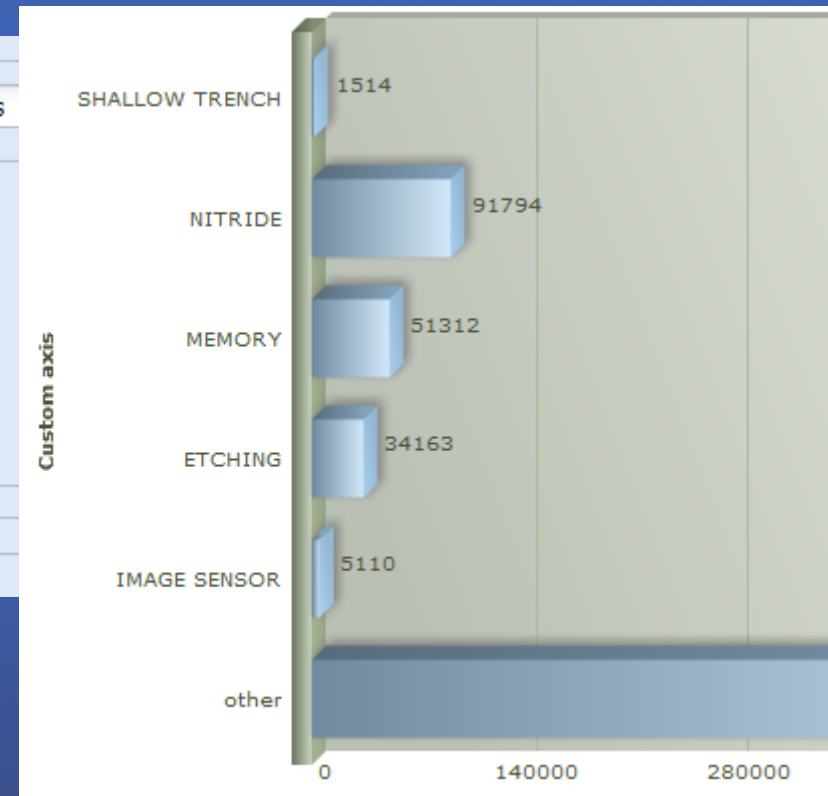
**Analysis axis**

Axis value: Custom      Axis name: Custom axis

**Custom Axis**

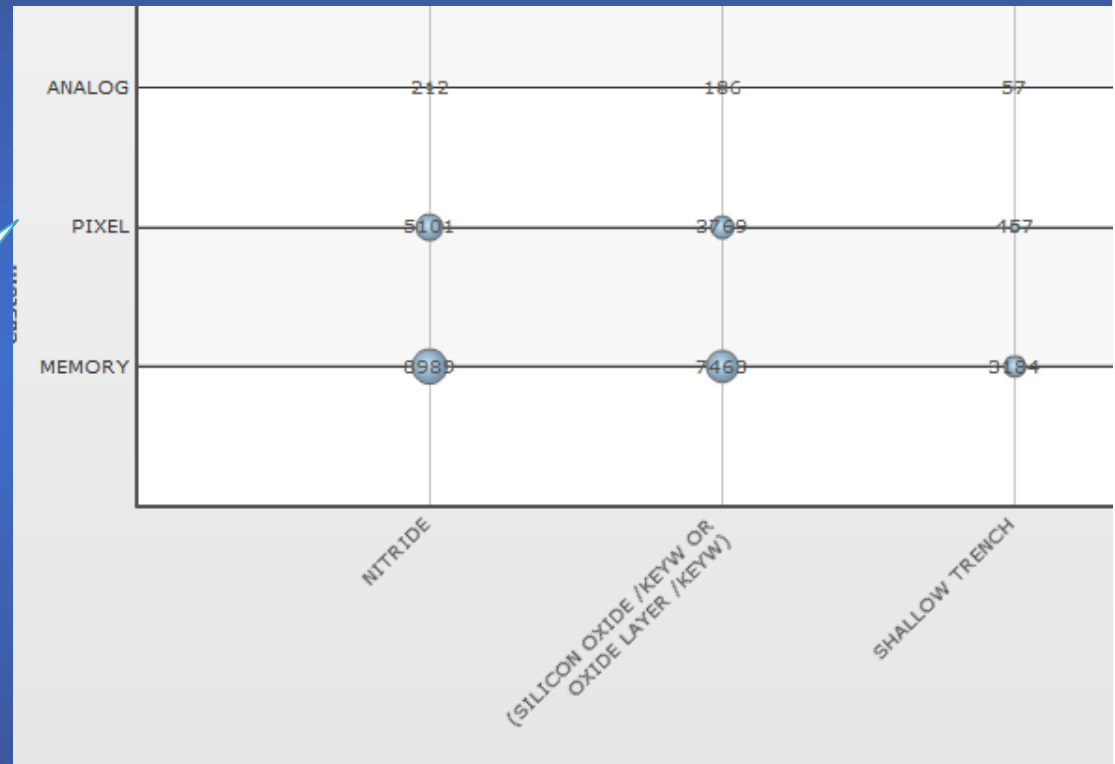
Custom	(shallow trench)/ICLM
Custom	nitride/Keyw
Custom	memory/iclm
Custom	etching/iclm
Custom	image sensor/iclm

Filter



# Crossing custom axis

- Crossing independent claim searches with concepts



Independent  
Claim  
content

Concept words

# Comparing scale-out and small corpus techniques

	Traditional /small corpus	Large scale
Pros	<ul style="list-style-type: none"><li>• Ability to edit the data</li><li>• Save edited data</li><li>• Maps</li></ul>	<ul style="list-style-type: none"><li>• Practically unlimited in volume</li><li>• Full integration with search engine</li><li>• Full access to text mining features</li></ul>
Cons	<ul style="list-style-type: none"><li>• Limited to 60,000 patent families</li><li>• Full text queries often limited</li><li>• Have to wait for data to be extracted/saved</li></ul>	<ul style="list-style-type: none"><li>• No data vetting / editing</li><li>• No data set saving</li></ul>

We need to get the best of both worlds

# Converged tool: orbit analysis module

- 3 data source
  - Live from database
  - Workfiles
  - Saved analysis: edited data
- Same features, same user interface

# Leveraging concepts: Mapping module



# Questions

- Feel free to share your horror / success stories in analyzing hundreds of thousand of patents.



Thank you very much