

1

Hybrid Search: Effectively Combining Keywords and Semantic Search

**Ravish Bhagdev, Sam Chapman,
Fabio Ciravegna, Vitaveska Lanfranchi**

web Intelligence technology lab,
nlp group, department of computer science,
university of sheffield

Daniela Petrelli

department of Information studies,
university of sheffield



Outline of Talk

- Aim of paper
 - Search for what and in what conditions? 3 mins+3 slides
- Hybrid Search as a way to overcome limitation of classic semantic search 7min+9slides
- Implementing Hybrid Search into K-Se 4 mins+5slides
- Experimental Evaluation in vitro and in 8 mins+15 slides
- Conclusion and future work 3 mins+2 slides

- We propose a search method
 - Designed for the Semantic Web
 - Seen as a collection of both documents and metadata,
 - Designed to achieve two tasks:
 - Document retrieval: searching for documents using concepts or keywords of interest
 - Knowledge retrieval: retrieving facts from a knowledge base (i.e. triples)



- Differently from [1, 2, 3, 4, 9],
 - We expect metadata to cover only partially the user information needs
 - Reasons:
 - limitations in the ontology wrt user needs
 - limitations in the annotation capabilities
 - i.e. limitations in IE capabilities
 - metadata unavailable for a specific document



- Semantic search as metadata-based search defined according to an ontology,
 - Annotations are unambiguous
 - OS Does not suffer from ambiguity and synonym issues of keyword-based systems (KS)
- Issue:
 - OS can fail to encompass user information needs
 - When metadata does not completely cover user needs



- We propose a model of searching combining
 - the flexibility of keyword-based retrieval
 - querying and reasoning capabilities of semantic search
- HS is formally defined as:
 - the application of semantic (metadata-based) search for the parts of the user queries
 - where metadata is available
 - the application of keyword-based search for the parts not covered by metadata.

- But also it must leave freedom to users to chose among the two paradigms!
 - As we will see users make a creative use of it



Queries in Hybrid Search

- Any boolean combination of three types of conditions

- pure semantic:

- via unique identification of objects/relations
 - e.g. via URIs or unique identifiers

- keyword-based

- matching on the whole document

- keyword-in-context

- matching keywords only within portion of documents semantically annotated with a specific type or instance

differently from other approaches (e.g. [9]), in HS conditions on metadata and keywords coexist.

Queries in Hybrid Search

- Any boolean combination of three types of conditions

- pure semantic:

- via unique identification of objects/relations

- e.g. via URIs or unique identifiers

- keyword-based

- matching on the whole document

- keyword-in-context

- e.g. it enables searching for the string "fuel" but only in the context of all the text portions annotated with the concept affected-engine-part [14]

differently from other approaches (e.g. [9]), in HS conditions on metadata and keywords coexist.

Example of Hybrid Query

$\forall x,y,z /$

(discoloration y) & (located-on y x) & (component x)

Querying Metadata

& (provenance-text-contains x "blade")

Keyword in Context Query

& (contains z "trailing edge") & (document z) & (provenance x z)

Keyword-based Query



Implementing HS: Indexing

- Documents are indexed using a standard keyword-based engine such as SolR
- Facts (e.g. extracted by an IE system) are stored in a Knowledge Base
 - e.g. a triple store like Sesame2 in the form of RDF triples.
- Provenance of facts recorded
 - E.g. As triples connecting
 - the facts' URIs and those of the document of origin
 - the facts' URIs and the original strings used in the documents



- Query is parsed and the different components (keywords, keywords-in-context and metadata) identified
 - keyword matches → traditional information retrieval system
 - metadata searches
 - Translated into a query language like SPARQL
 - Sent to a triple store
 - keywords-in-context queries
 - matched with provenance of annotations in documents
 - E.g. Using SPARQL and a triple store
- Finally, results are merged, ranked and displayed



- Merging keyword and semantic results is not straightforward
 - Keyword matching returns an ordered set of URIs of documents
 - a semantic search returns an unordered set of assertions < subj, rel, obj >
- Merging is a different task if:
 - Document Searching
 - Returns documents
 - Knowledge Searching
 - Returns triples



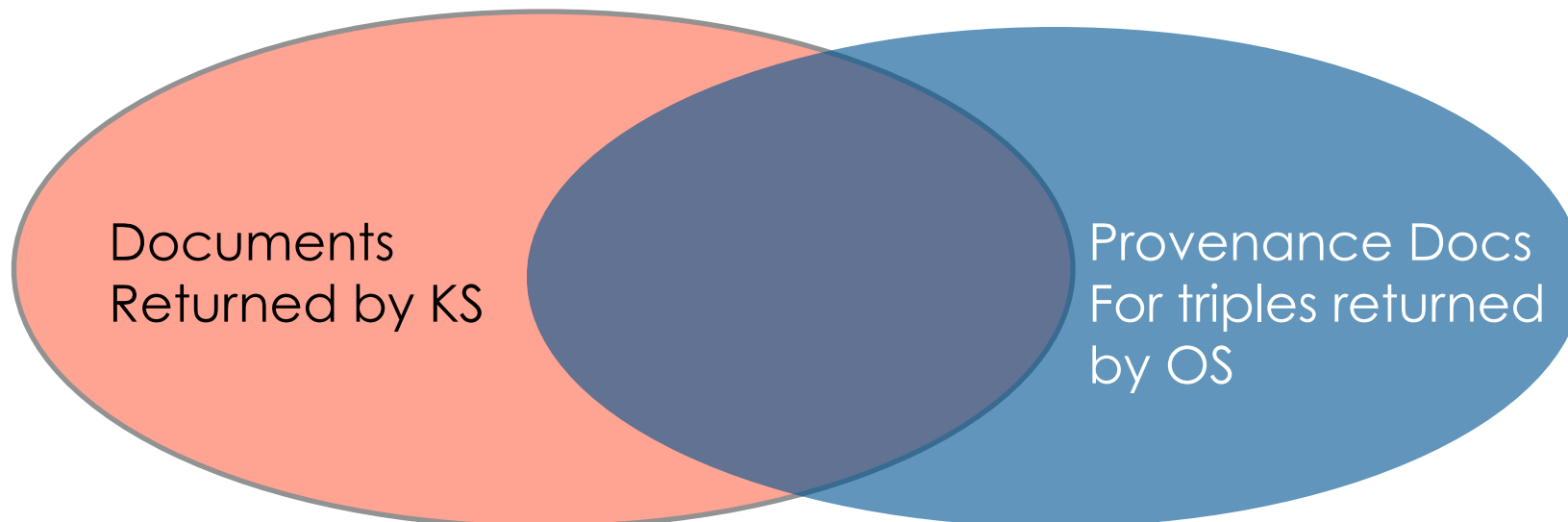
- Provenance of triples returns document ids for triples (URIs)

- Document Searching:

- Provenance URI set is intersected with URIs of documents returned by keywords

- $\text{HybridSearchUriSet} = \text{KSDocUriSet} \cap \text{OSDocUriSet}$

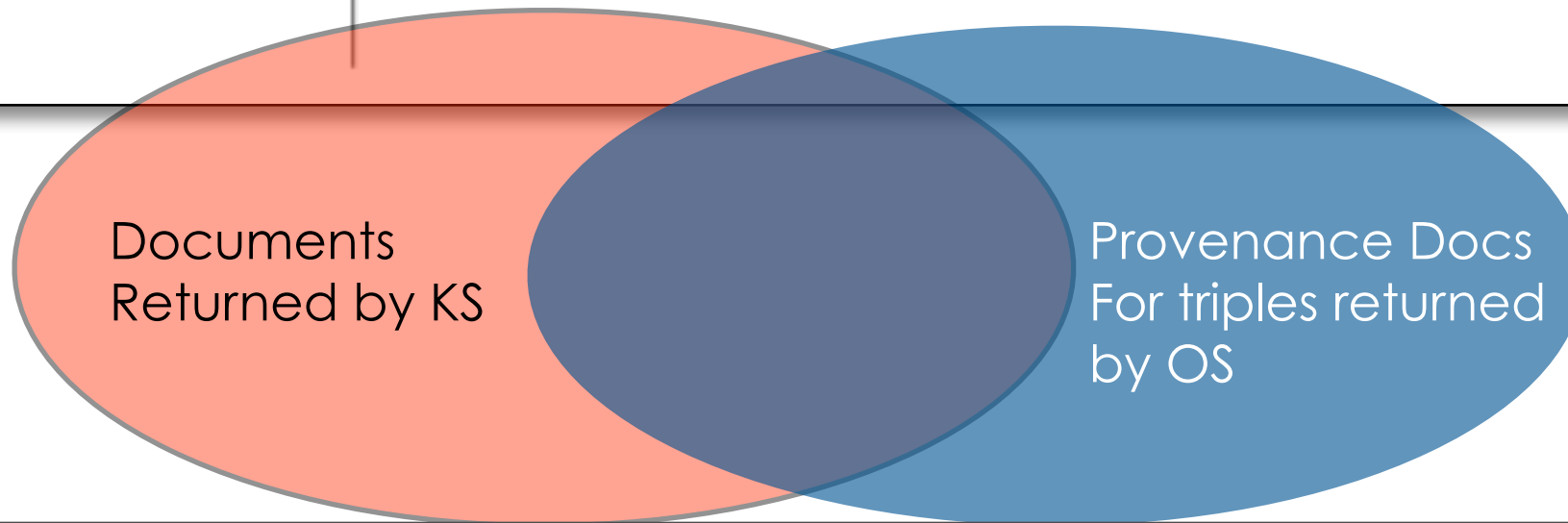
I won't mention ranking here



- Provenance of triples returns document ids for triples (URIs)
 - Knowledge Searching
 - Triples returned by semantic search are filtered to remove those whose provenance does not point to any of the documents returned by the keywords

I won't mention ranking here

```
HSTripleSet = All triples ∈ OSTripleSet  
              Where Provenance(triplei) ∈ KSDocUriSet
```



Expected effect of HS: Document Searching

14

- With respect to OS
 - Recall expected to increase
 - Use of keywords where metadata is missing enables to answer otherwise impossible queries
 - Precision may suffer because of polysemy
- With respect to KS
 - Precision and recall expected to increase
 - Ambiguity and synonymy are dealt with by semantic search when available
 - Higher recall and precision
 - As keywords are combined with metadata in the same query, the context given by the available metadata helps in disambiguating keywords as well
 - higher precision



Expected effect of HS : Knowledge Searching

15

■ With respect to OS

■ Precision increased

- Use of keywords where metadata is missing enables more precise queries
 - although less precise than the ideal ones

OR

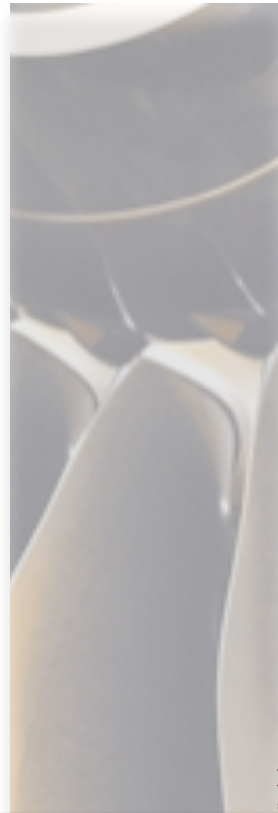
■ Recall increased

- Use of keywords where metadata is missing enables to answer otherwise impossible queries

■ Precision may suffer because of polysemy

■ With respect to KS

- KS does not cover Knowledge Searching

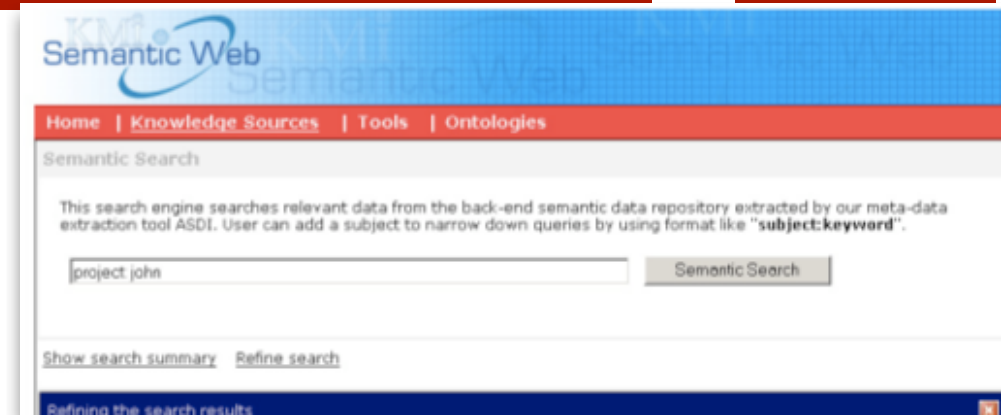


Next slide:
We have
implemented a
version to confirm
our expectation

Implementing HS: What Search Strategy?

- Keyword-based approaches

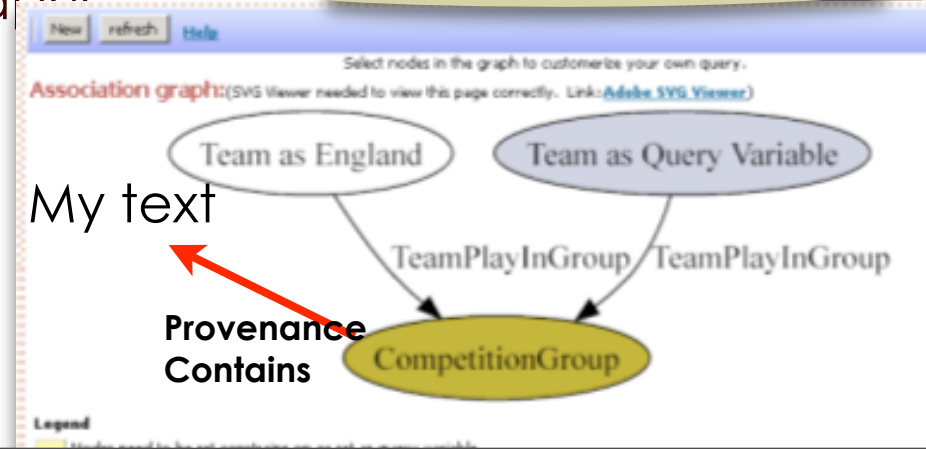
- Require translating all the keywords in order to perform the query
 - E.g. SemSearch
 - HS implemented by replacing keywords in the query with classes from the ontology when possible while leaving the rest for pure keyword based searching
 - Keywords in context rather difficult



Go through this and next slide very quickly !!

- View-based approaches

- Based on querying by building visual graphs
 - E.g. Falcon
 - HS support by adding two arc types
 - document-contains
 - Object description contains



Search Strategy (ctd)

17

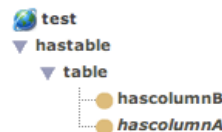
- A natural language approach
 - E.g. Aqua
 - HS supported by recognising expressions like
 - “and the document contains...”
 - And its description contains
- Form-based approaches
 - HS supported by introducing
 - Keyword Search field
 - Enable keyword Matching on fields

The screenshot shows a web interface for 'Question Answering'. At the top, there is a search bar with the text 'Show me all planet stories written by a researcher in AKT' and an 'Ask!' button. Below the search bar, there is a checkbox for 'Make Use of Learning Mechanism for relations' which is checked. The main content area displays the results of a 'Relation Similarity Service' query. It shows a 'Linguistic Triple' and an 'Ontology Triple' with their respective components and relationships. The linguistic triple is 'planet stories - written - researcher - akt'. The ontology triple is 'kmi-planet-news-item - has-author owned-by - researcher - akt'. There are also two notes explaining the mapping of the learning mechanism to the ontology terms.

- Form-based implementation of hybrid search initially created for Jet Engine Designers

k·now

Available Reports



The screenshot displays a web browser window with the IPAS (Incident Processing and Analysis System) interface. The browser's address bar is empty, and the window title is partially visible as "The University Of Sheffield".

The interface is divided into two main sections:

- Search Section (Right):** Features a "Search" tab and a "Graph" tab. The "Search" tab contains a "Keyword Search:" input field with the placeholder text "(optional)". Below this is a "Number of results per page" dropdown menu set to "10". A "SEARCH" button is located at the bottom of this section. A note below the search options reads: "[Click on an ontology concept (left) to add search criteria]".
- Navigation Tree (Left):** Titled "Event Report", it lists various ontology concepts for selection. The tree structure is as follows:
 - Event Report
 - Report Number
 - Report Creation Date
 - Report Author
 - Referred Service Event
 - Service Event
 - Event Date
 - Event Type
 - Event Category
 - Operational Effect
 - Flight Regime
 - Event Location
 - Airframe Cycles
 - Airframe Hours
 - Engine Installed Location
 - Fuel Dumped
 - Delayed Time
 - Affected Engine
 - Engine
 - Engine Serial Number
 - Engine Type
 - Installed Part
 - Component
 - Removed Part
 - Event Description
 - File Location

K-Search

18

The screenshot displays the IPAS (Integrated Platform for Aircraft Systems) web application. The interface includes a navigation menu on the left, a search bar at the top, and a main content area. The search query is "FMU OR fuel flow transmit", and the results are visualized in a 3D pie chart titled "Analysed Graph".

Search Query: FMU OR fuel flow transmit

Analysed Graph Data:

| Component | Percentage |
|-----------------------------------|------------|
| fuel metering unit (fmu) | 41% |
| fuel flow transmitter | 33% |
| fuel metering unit (fmu) | 14% |
| fuel flow transmitter (sb73-c579) | 5% |
| fmu | 7% |



- We have performed 2 types of evaluations using K-Search:
 - in vitro:
 - Effectiveness of query strategy with respect to standard KS and OS
 - in vivo: testing the system with real users
 - 32 users Rolls-Royce engineers
 - Evaluation enables verifying suitability for use in a real environment

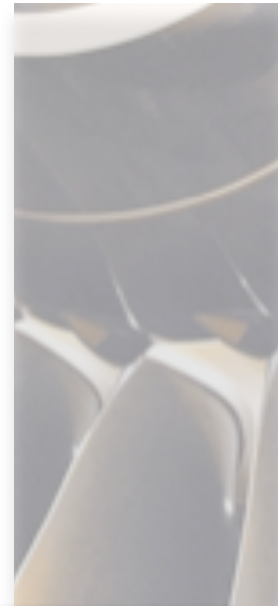


- Automatic extraction of information from event report
 - 18,000 documents analysed
 - Mainly Forms implemented in Word
- Metadata generated according to an ontology developed by Aberdeen U
- Automatic extraction of metadata and indexing of documents

IE unable to cover all the ontology with sufficient accuracy

Applying information extraction

- AktiveMedia to annotate texts
- TRex system (Jiria et al. 2006) to train and extract
 - <http://tyne.shef.ac.uk/t-rex/>
- IE captures all the information in tables
 - 99% of the information captured (recall=99)
 - 98% of proposed information is correct (precision=98)



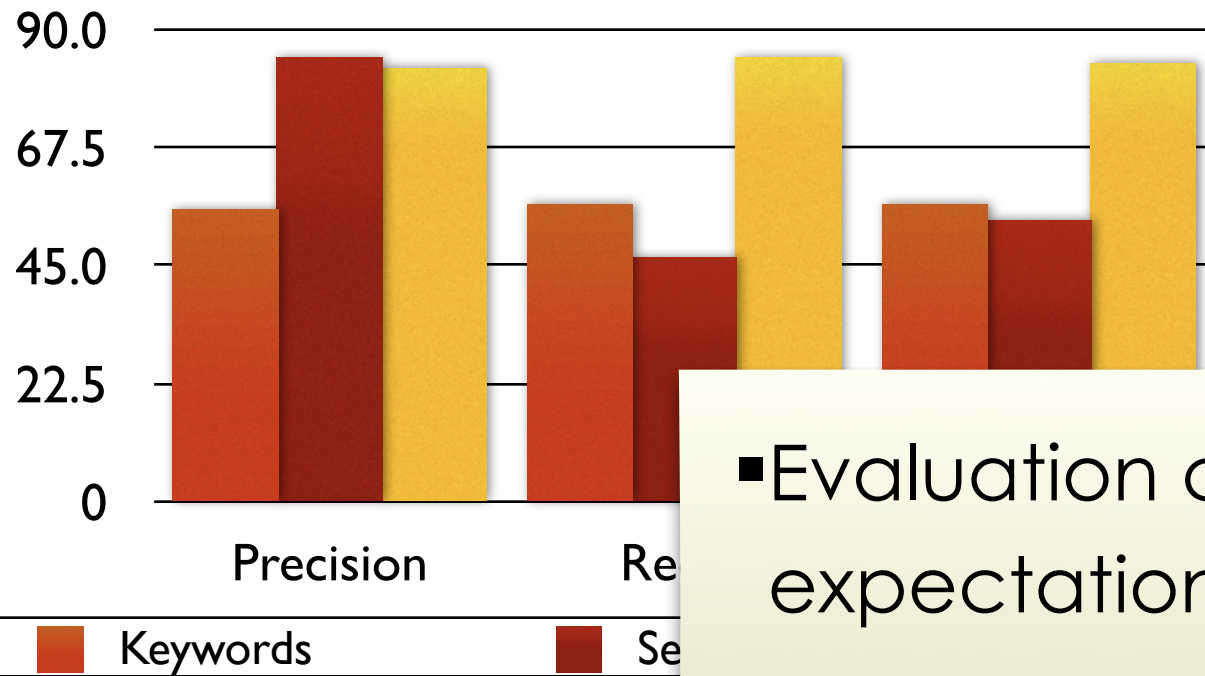
| | POS | ACT | CORR | WRONG | MISSED | PREC | REC | F1 |
|--------------------------------|-------------|-------------|-------------|-----------|-----------|-----------|-----------|-----------|
| airport | 120 | 120 | 120 | 0 | 0 | 100 | 100 | 100 |
| has_airframe_cycles | 104 | 104 | 104 | 0 | 0 | 100 | 100 | 100 |
| has_airframe_hours | 104 | 104 | 104 | 0 | 0 | 100 | 100 | 100 |
| has_author | 120 | 120 | 120 | 0 | 0 | 100 | 100 | 100 |
| has_engine_serial_number | 120 | 120 | 120 | 0 | 0 | 100 | 100 | 100 |
| has_engine_type | 120 | 120 | 120 | 0 | 0 | 100 | 100 | 100 |
| has_event_date | 120 | 120 | 120 | 0 | 0 | 100 | 100 | 100 |
| has_event_report_no | 356 | 358 | 356 | 2 | 0 | 99 | 100 | 100 |
| has_part_description_installed | 120 | 113 | 111 | 2 | 9 | 98 | 93 | 95 |
| has_part_description_removed | 120 | 133 | 120 | 13 | 0 | 90 | 100 | 95 |
| has_part_number_installed | 120 | 113 | 111 | 2 | 9 | 98 | 93 | 95 |
| has_part_number_removed | 120 | 133 | 119 | 14 | 1 | 89 | 99 | 94 |
| TOTAL | 1644 | 1658 | 1625 | 33 | 19 | 98 | 99 | 98 |

- 21 topics of search, discussed with users, e.g.
 - "How many events were caused during maintenance in 2003?"
 - "What events were caused during maintenance in 2003 due to control units?"
 - 'Find all the events associated with damage to acoustic liners following bird strike'
- Queries:
 - "what events caused during maintenance in 2003 were due to control units?"
- Translated into a set of queries in KS, OS and HS



K-Search on Event Reports

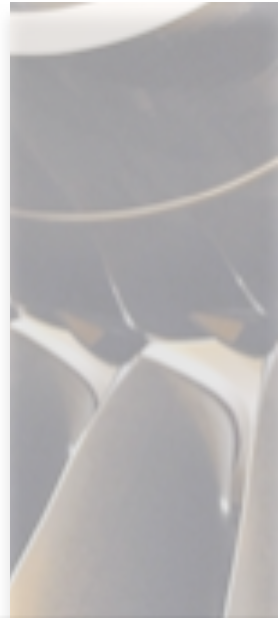
- Accuracy in the first 20 hits on a sample of 400 docs



- Evaluation confirms our expectation:

- Similar results for 50 hits

- Higher recall wrt OS and KS
- Higher precision wrt KS
- Slightly lower precision wrt OS



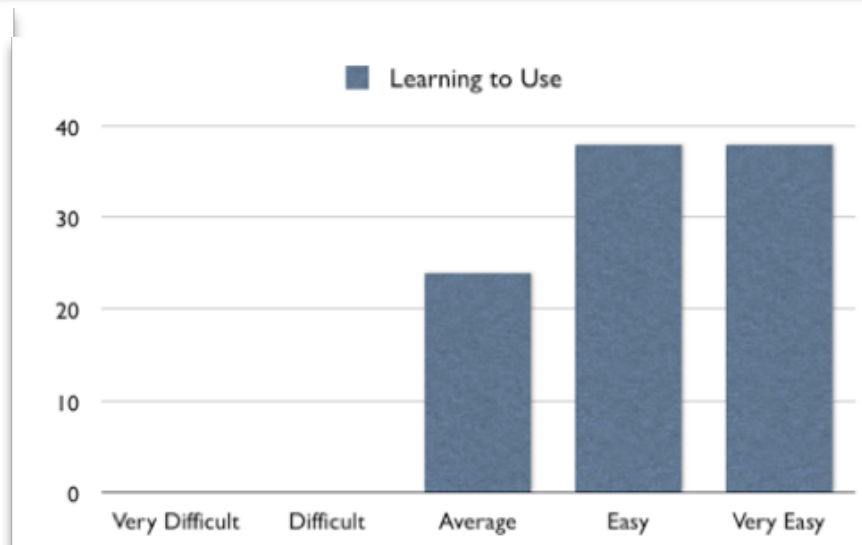
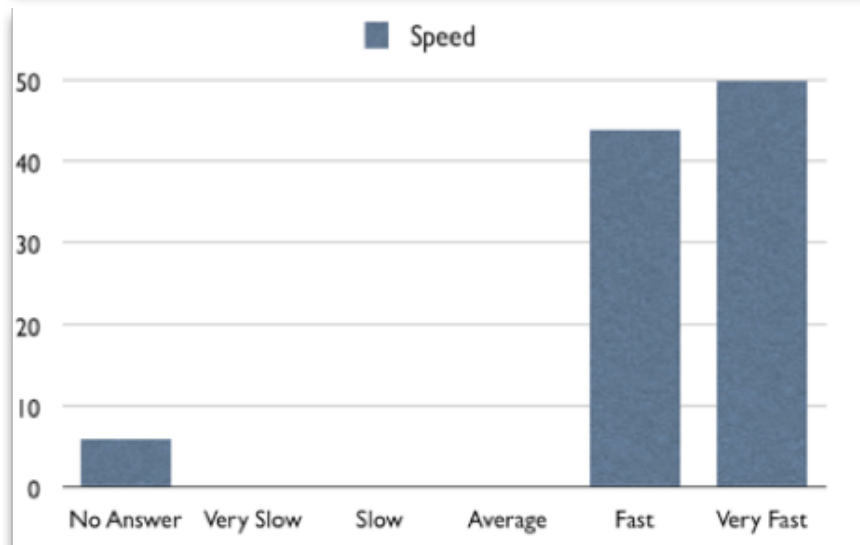
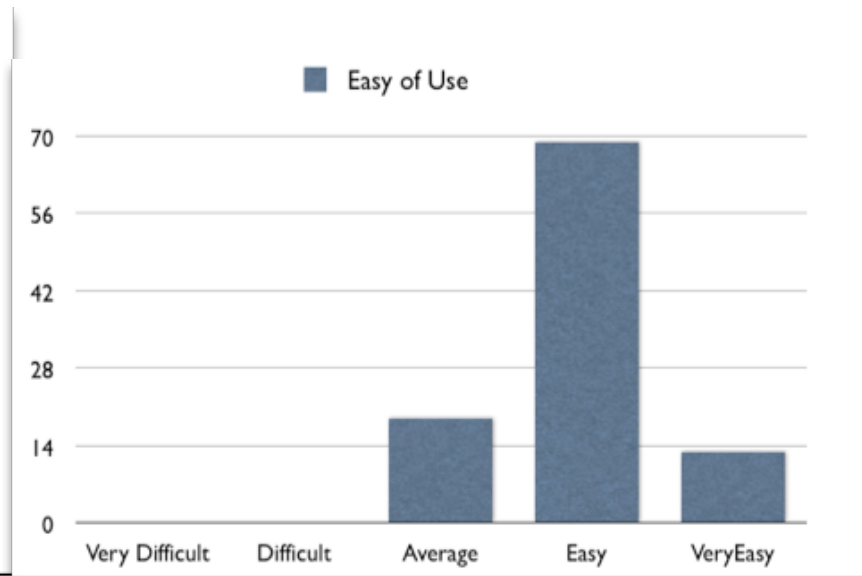
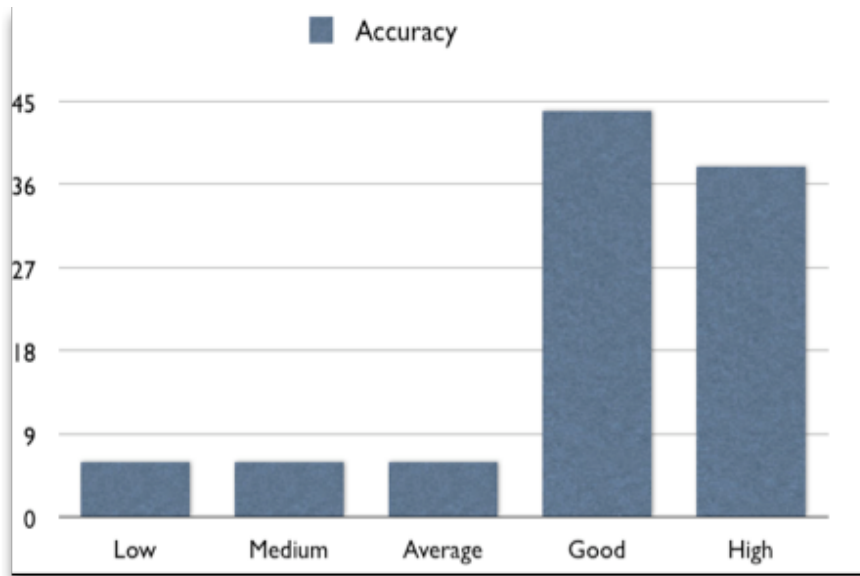
- Goal: verifying suitability for use in a real environment
 - 32 users Rolls-Royce engineers from different parts of the company
 - 90 minutes of test
 - Short introduction
 - 3 monitored tasks
 - One given (including solution)
 - One given (no solution)
 - One free task
 - Availability of system on intranet for the following period
- Evaluation: video recording, interview + log analysis



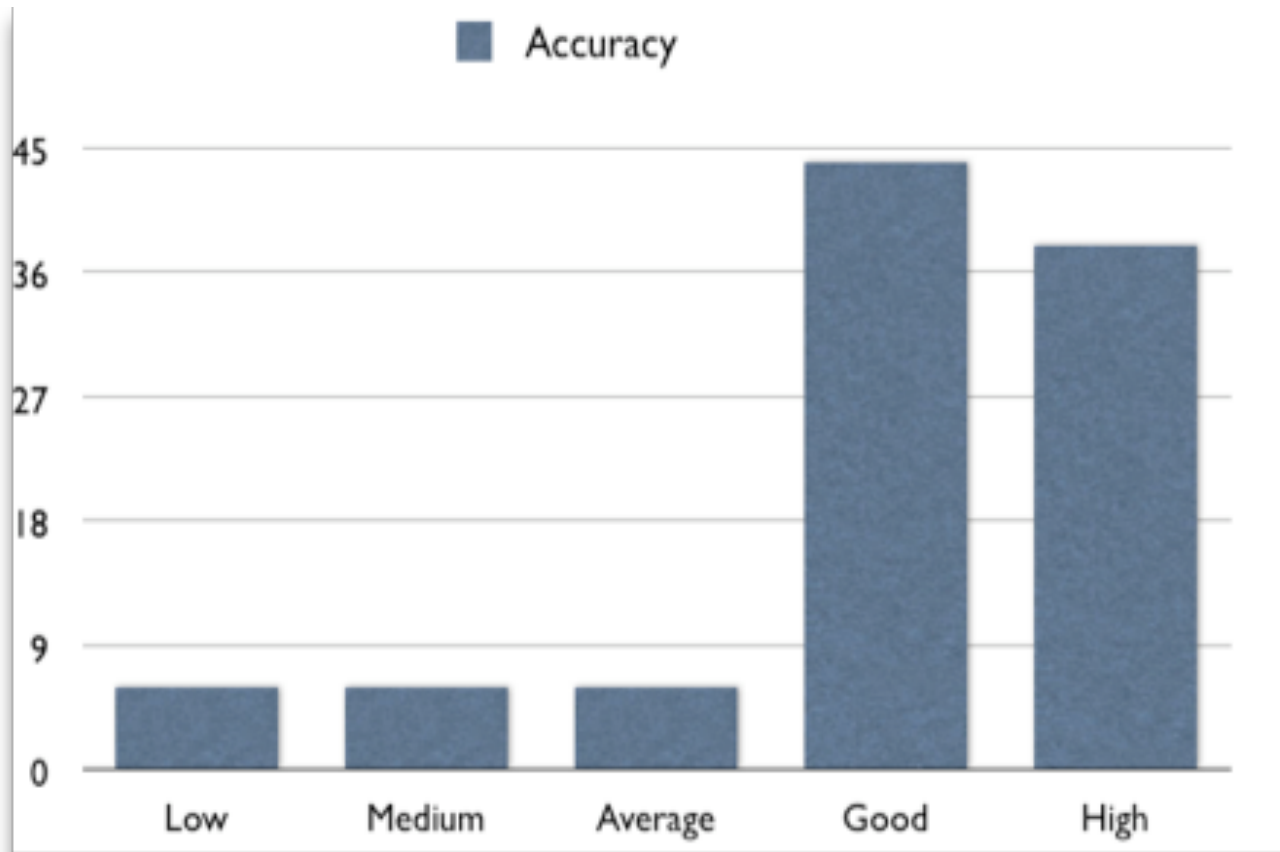
- Do user understand the hybrid paradigm?
- Are they able to search using HS?
- Do they actually use HS when confronted with a real searching task?
- Would the users be willing to use the system for their everyday work?



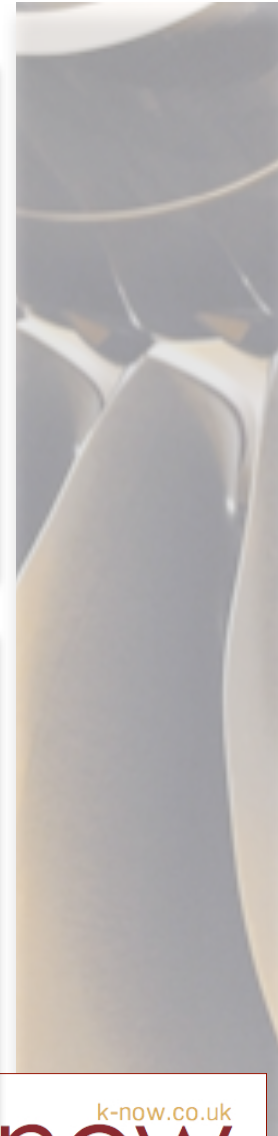
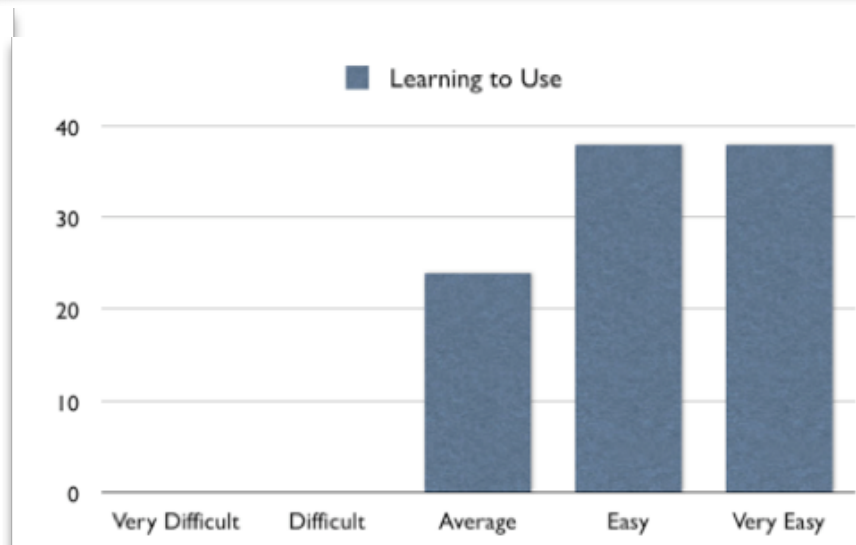
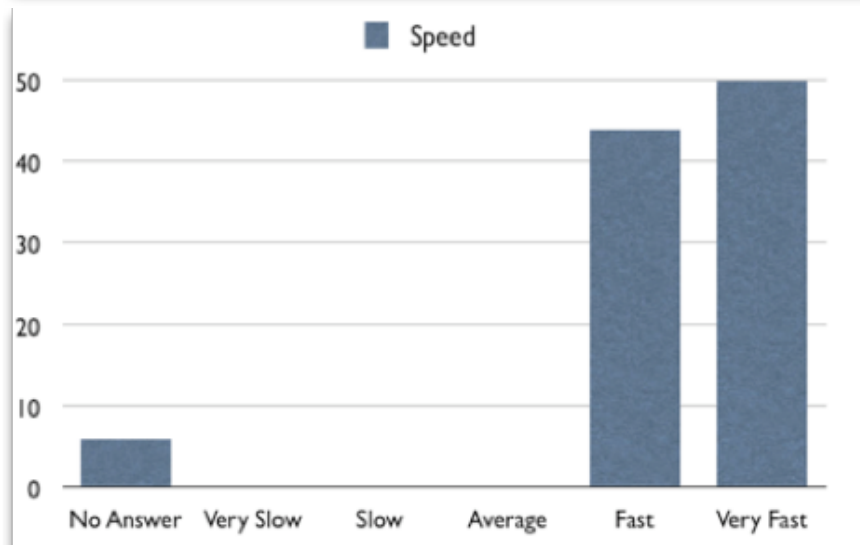
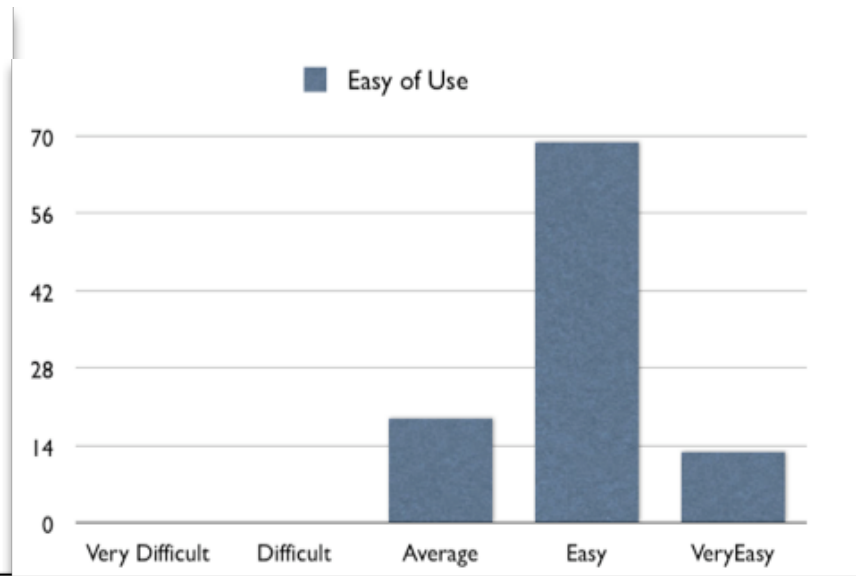
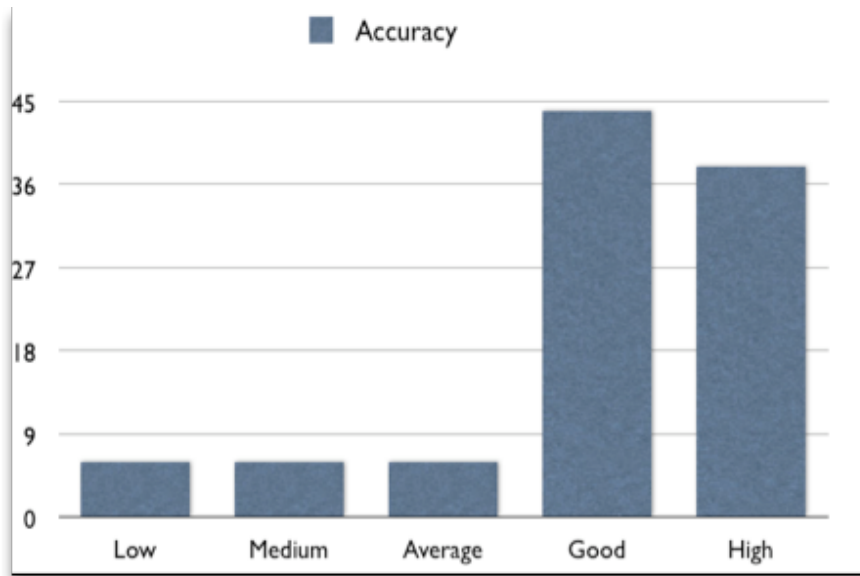
Liked by the users?



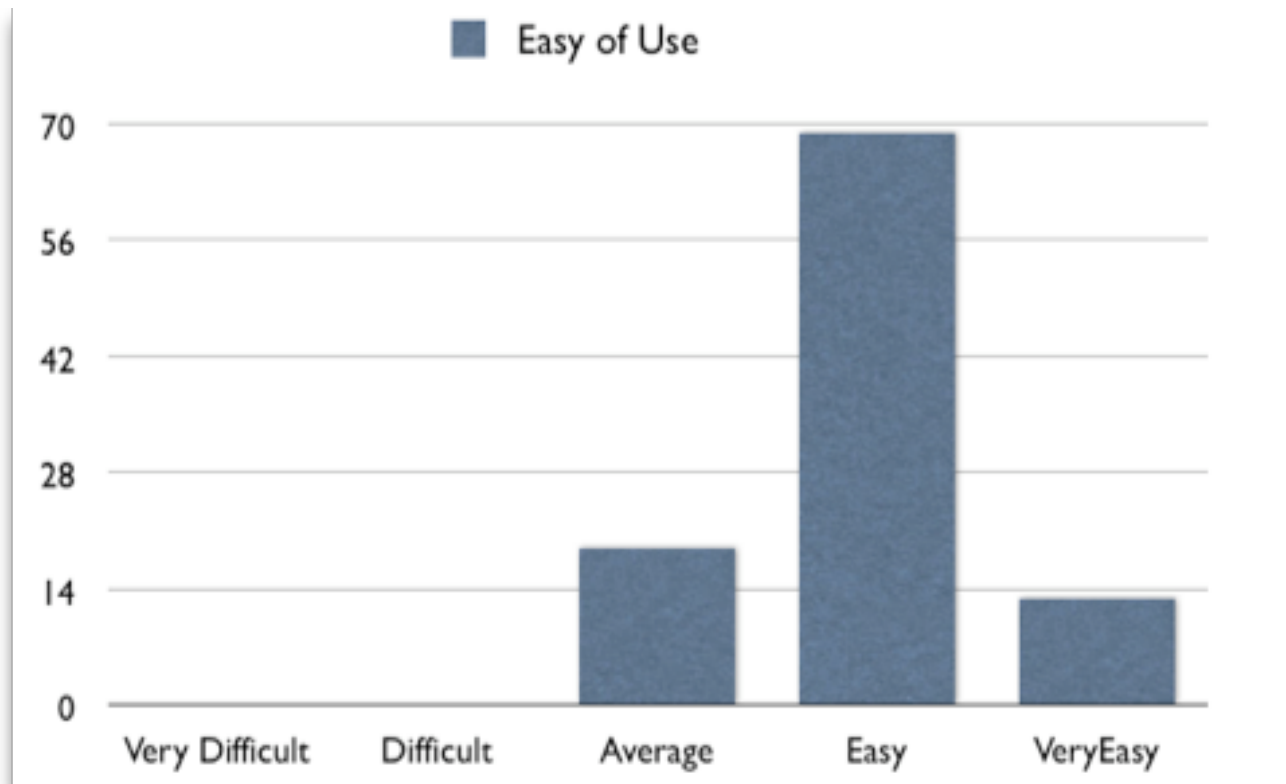
Liked by the users?



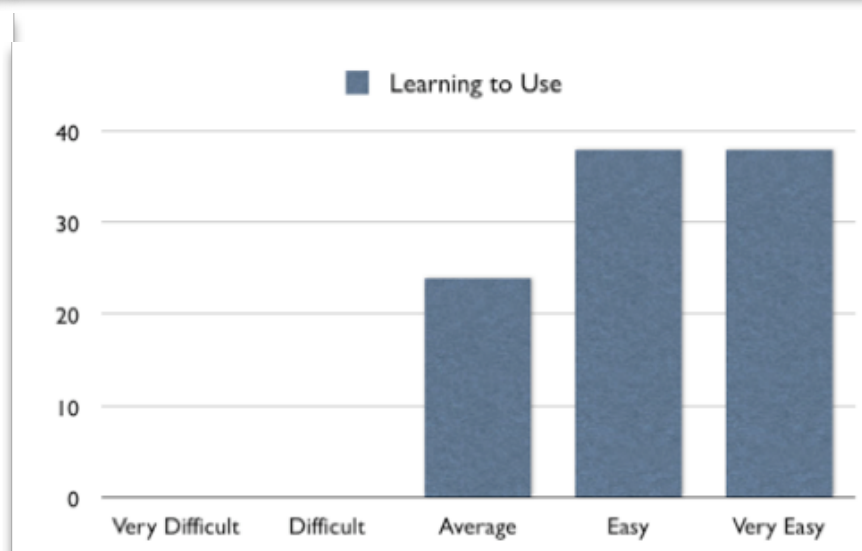
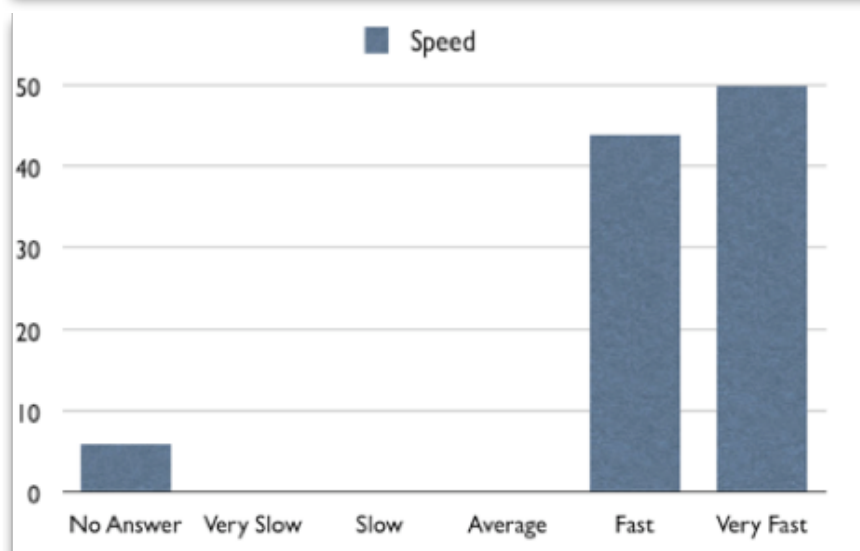
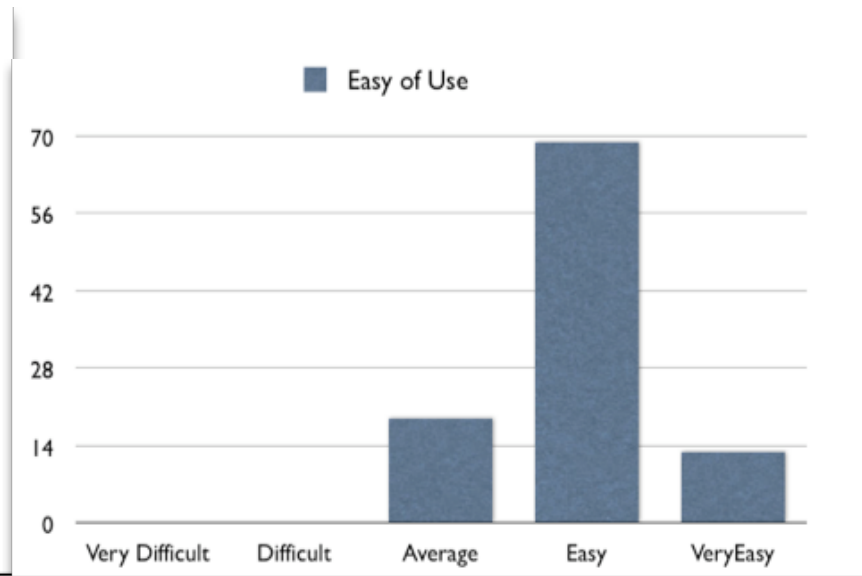
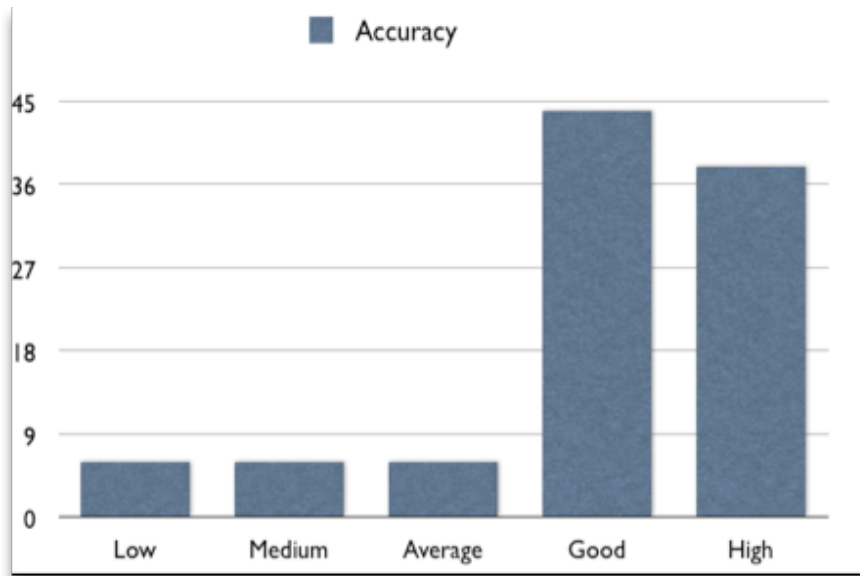
Liked by the users?



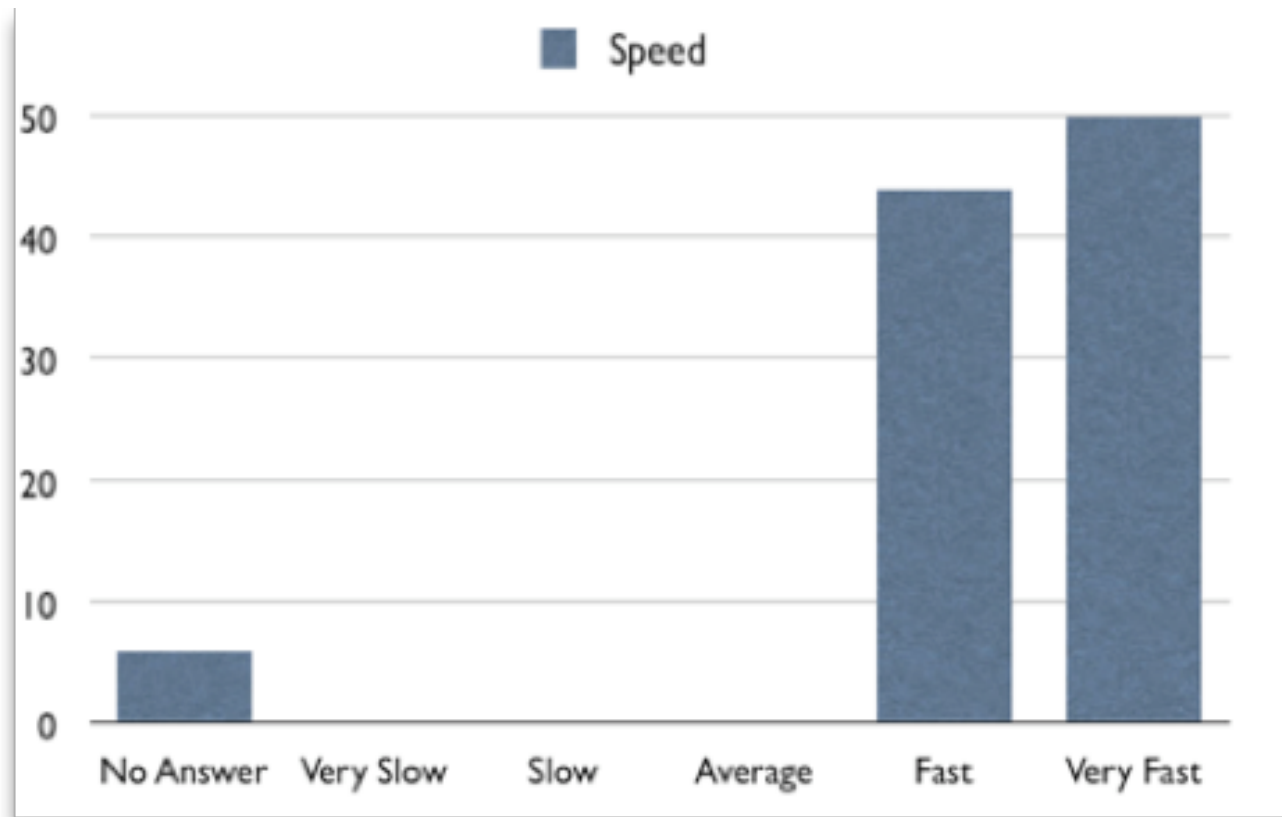
Liked by the users?



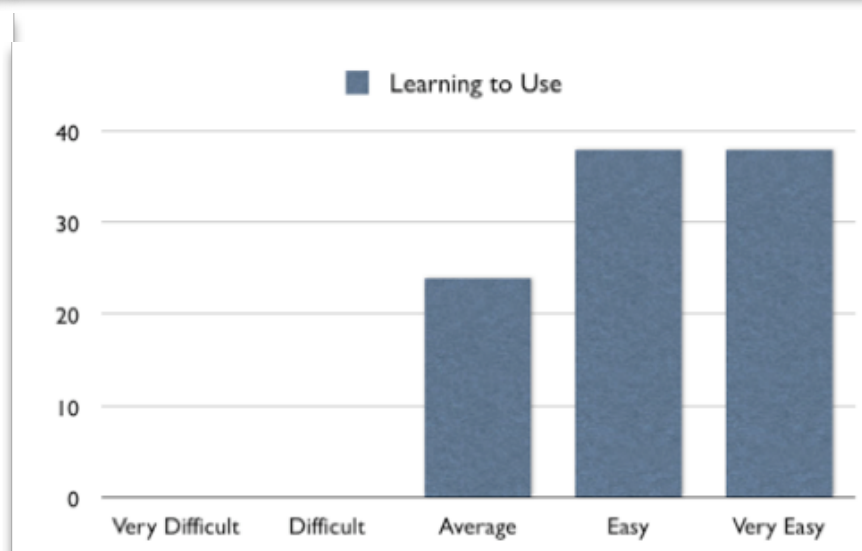
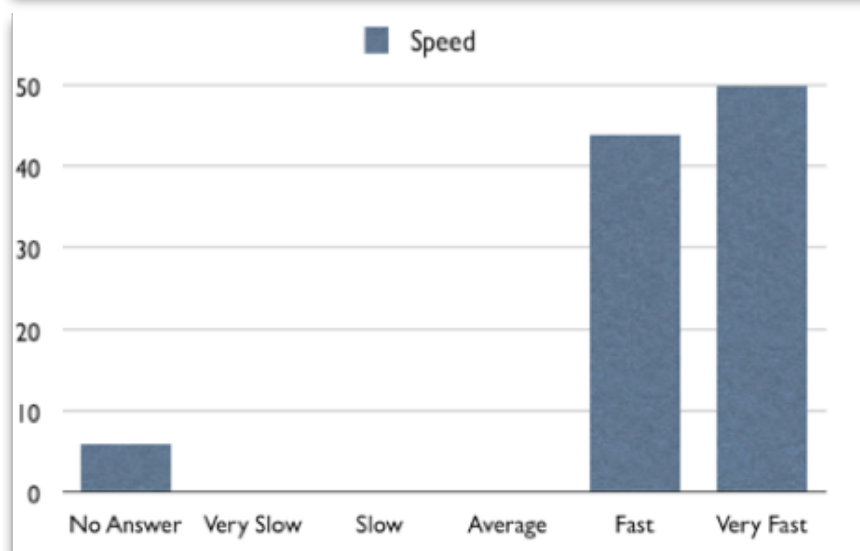
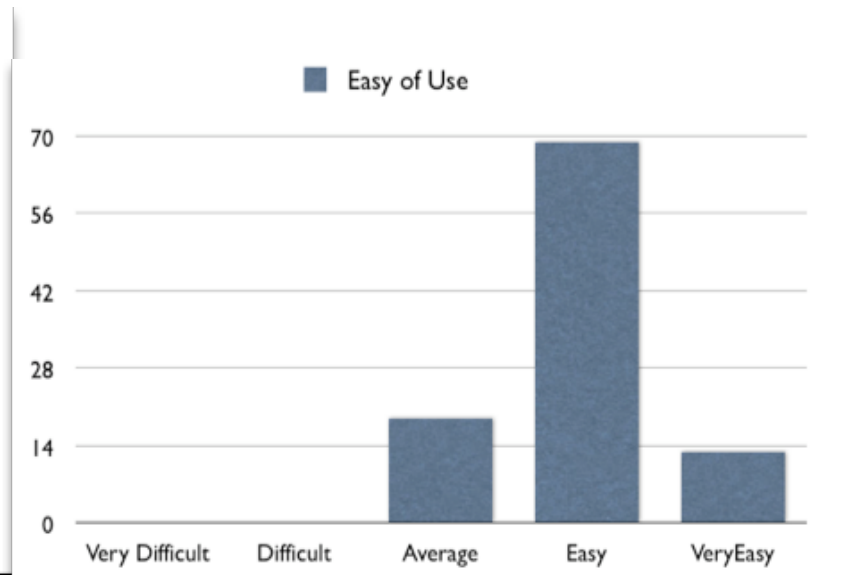
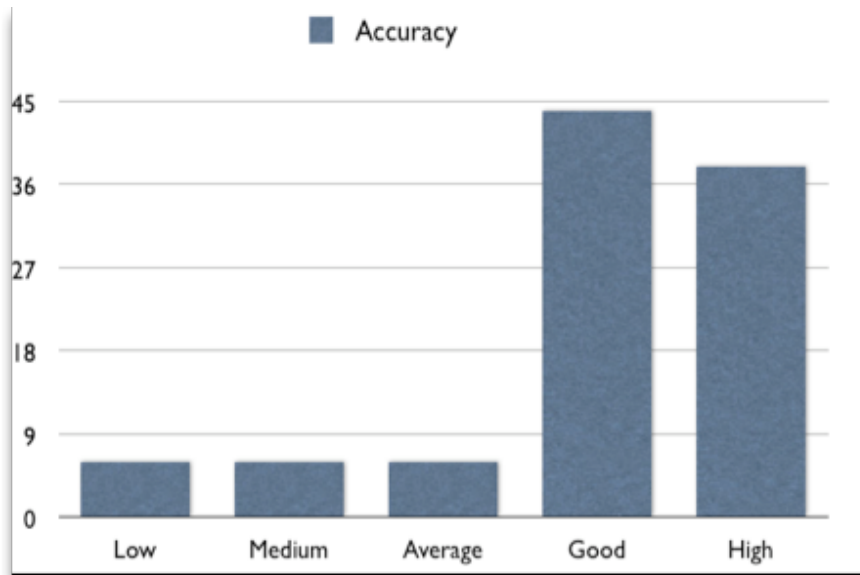
Liked by the users?



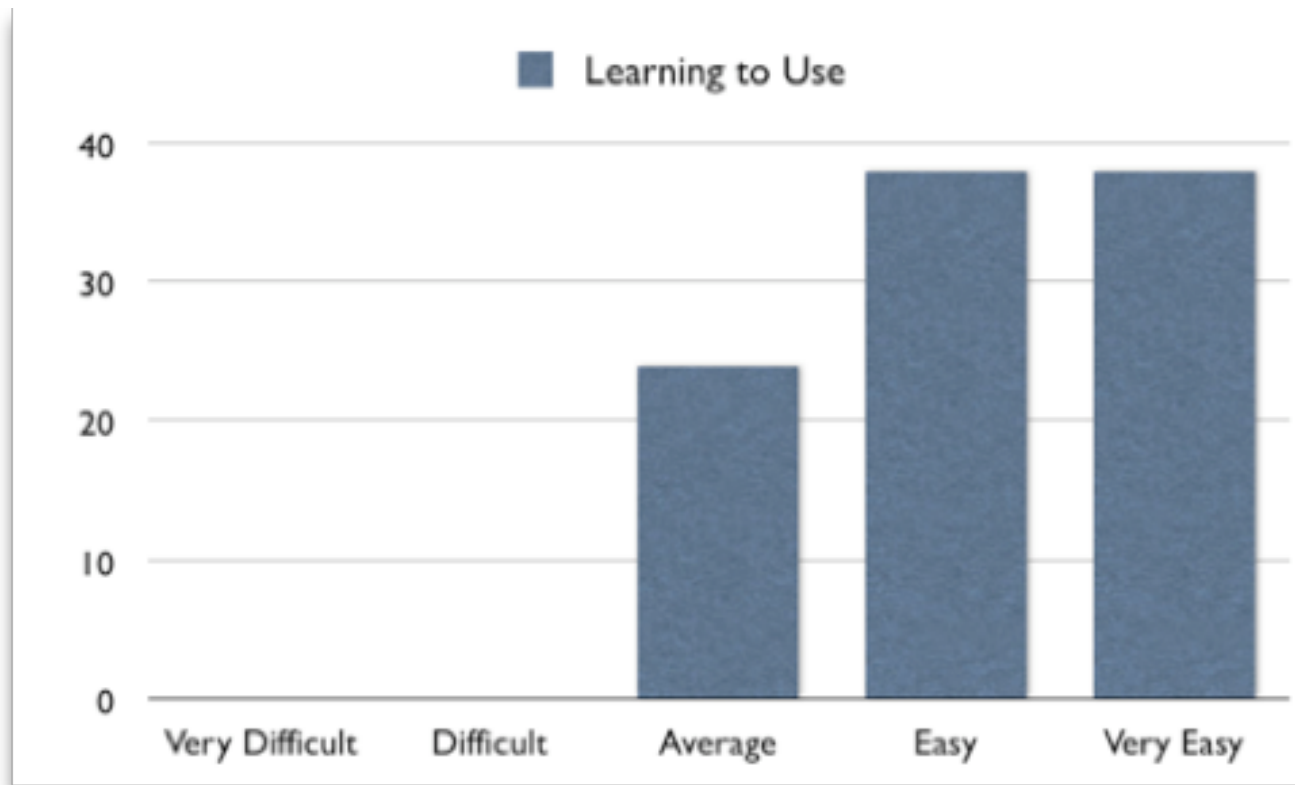
Liked by the users?



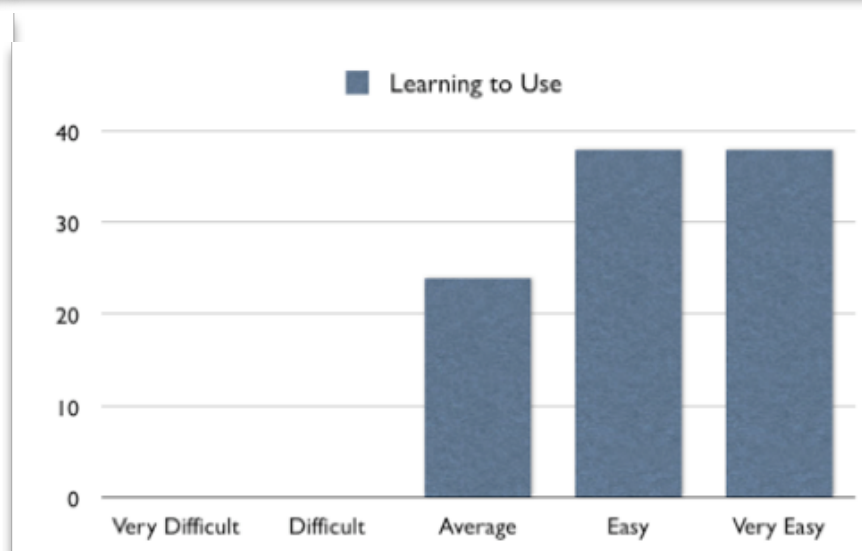
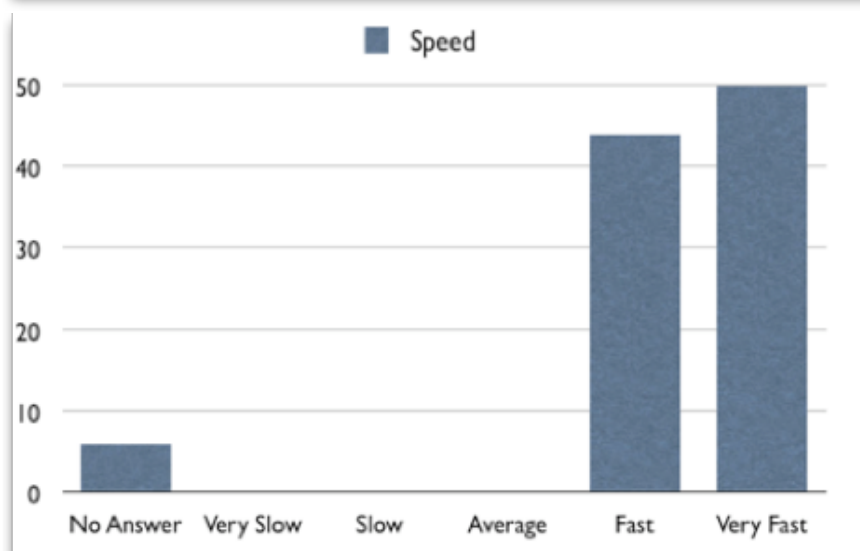
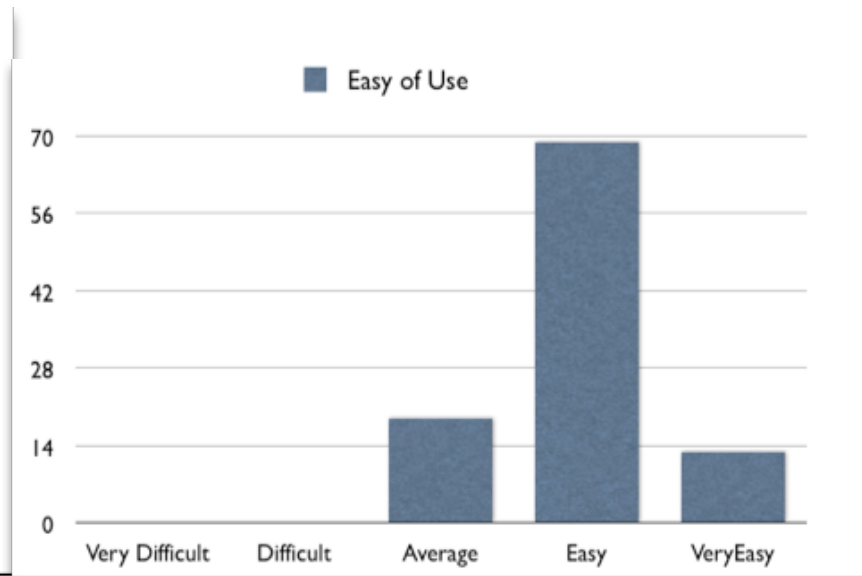
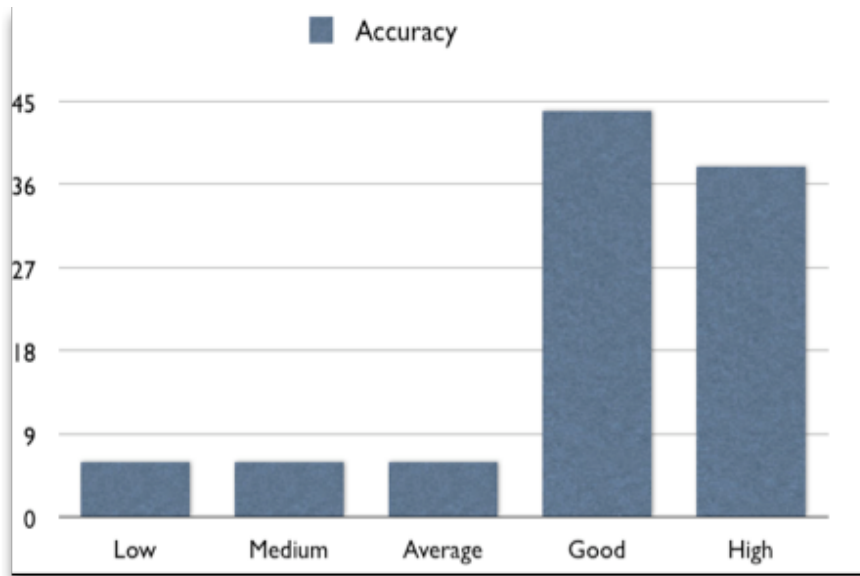
Liked by the users?



Liked by the users?



Liked by the users?

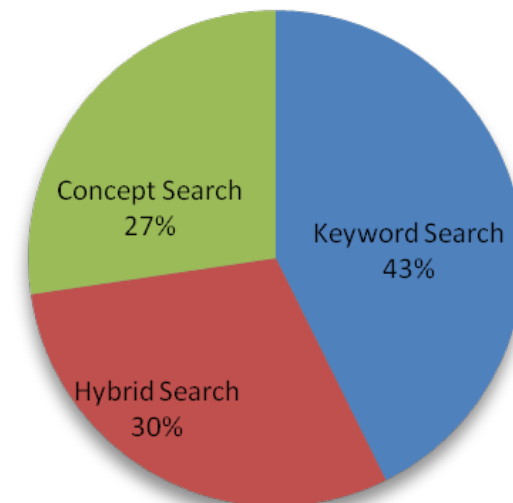


- Service engineers showed a clear predilection for hybrid search:
 - 61% of the search were executed using the hybrid modality
 - 24% using semantic search
 - 15% using keyword search.

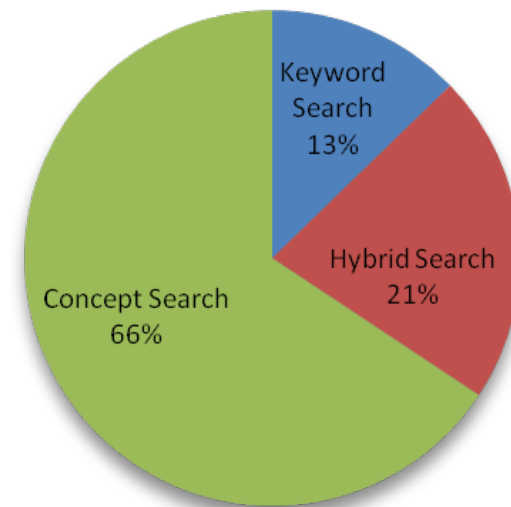
Reason: data they were looking for was not all covered by the metadata



- Designers tended instead to favour keyword search:
 - 43% of the searches were executed using keyword search
 - 30% using hybrid
 - 27% using semantic search.



- The users belonging to other groups showed a predilection for concept search:
 - 66% of the searches were executed using semantic search
 - 24% using hybrid
 - 15% using keyword search.



- Finalist of Rolls-Royce Director's Creativity Award 2007
 - Voted by employees for its innovation potential



Liked by Users?

- Support to the design of new Trent XWB
 - Porting to 9 Information Sources
 - 2008-2009
 - Carried out by:
 - 50% University
 - 50% k-now ltd (university spinout-company)
- Funds requested to UK Government for use of K-Tools for use in manufacturing



■ Hybrid Search

- It is compatible with the most used semantic search paradigms
 - Overcomes limitation of most current approaches based on metadata only
- Accommodates different search strategies
 - Users can choose how to perform the query
- Experimentally definitely outperforming both KS and OS



- Search across linked ontologies over intranet
- New ways of capturing information
 - User centred for new data
 - Cross-media
 - K-Forms
 - IE for legacy data
 - Cross-media



Thank You!

Contact Information

- www.dcs.shef.ac.uk/~fabio
- fabio@dcs.shef.ac.uk

Intelligent Web Technologies Lab

- <http://nlp.shef.ac.uk/wig/>

NLP Sheffield

- <http://nlp.shef.ac.uk/>

University of Sheffield

- www.shef.ac.uk

K-Now Ltd

- www.k-now.co.uk

