

DB ✕ HCI: Towards Bridging the Chasm Between Graph Data Management & HCI

Sourav S Bhowmick

Nanyang Technological University
Singapore

www.ntu.edu.sg/home/assourav

assourav@ntu.edu.sg, sourav@mit.edu



The World has Changed

Then

- Data is generated in companies
- Resides in companies
- Used by companies
- DB-literate users

Now

- Data is generated by everyone
- Resides everywhere
- Used by everyone
- Non-DB literate users





The World Has Changed

“The old computing was about what computers could do; the new computing is about what people can do”

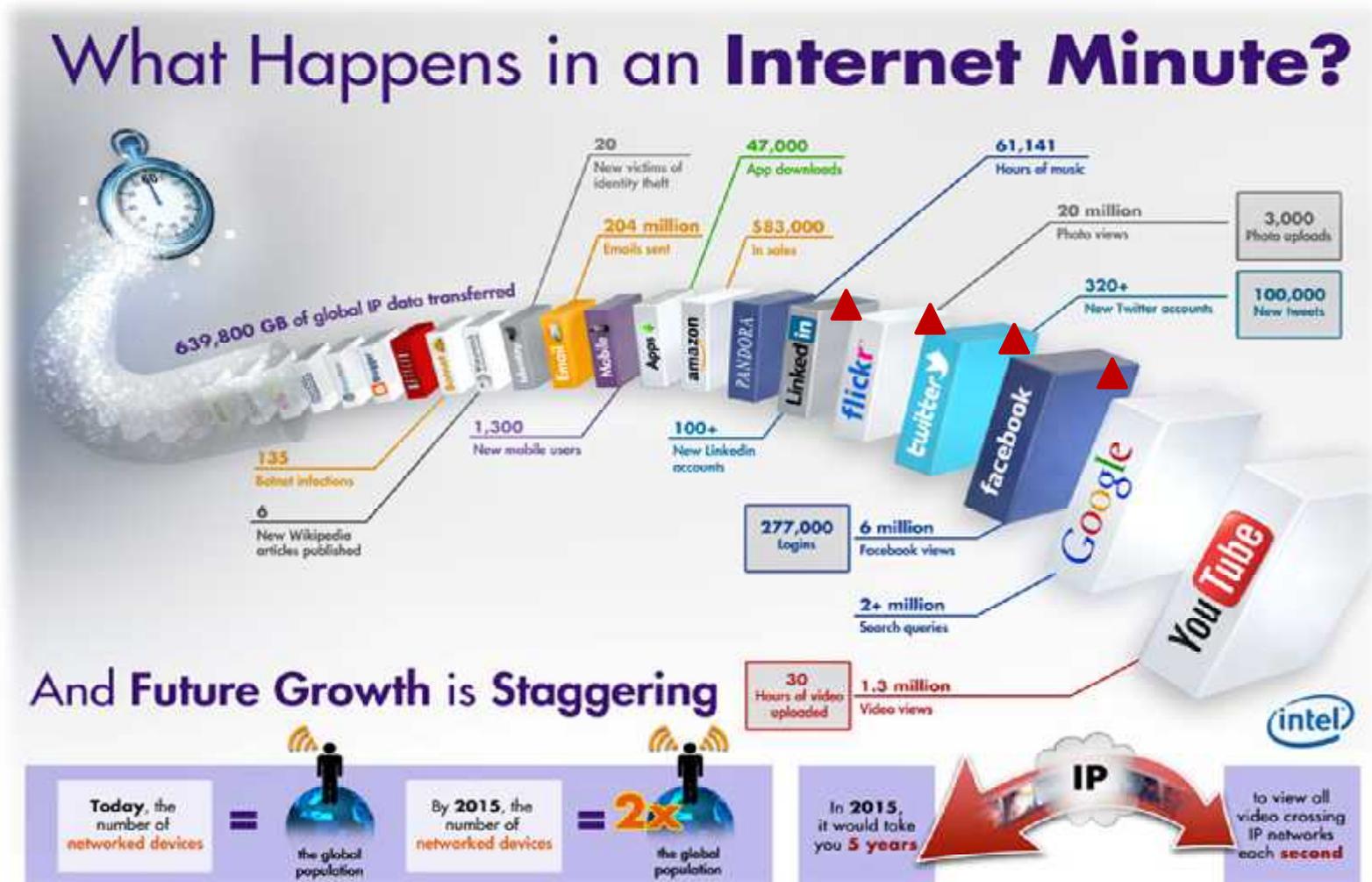
can do,, *Ben Shneiderman (2002)*

computing is about what people





Users are willing to share information online



Source: <http://www.intel.sg/content/www/us/en/communications/internet-minute-infographic.html>



The Road Ahead



- DB in the changing world
- HCI in the changing world
- The chasm!
- HCI-aware data management
- Conclusions



Complexity has Increased



Data streams



DB+ IR

mobile data management



Parallel & Distributed DB

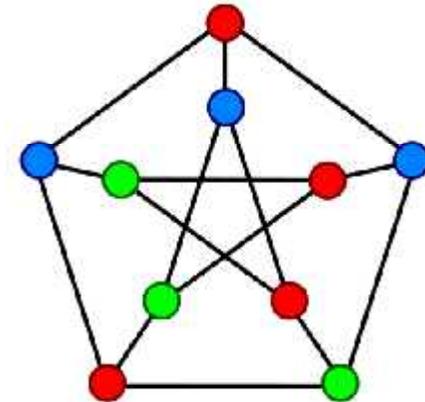


Cloud data management



<?xml?>

semistructured
data management



graph
data management

Object-Oriented Model

Object 1: Maintenance Report		Object 1 Instance	
Date		01-12-01	
Activity Code		24	
Route No.		1-95	
Daily Production		2.5	
Equipment Hours		6.0	
Labor Hours		6.0	

Object 2: Maintenance Activity			
Activity Code			
Activity Name			
Production Unit			
Average Daily Production Rate			

OO DB

Key	Product ID	Price (\$)	Prob.
a ₁	a	120	0.7
a ₂	a	80	0.3
b ₁	b	110	0.6
b ₂	b	90	0.4
c ₁	c	140	0.5
c ₂	c	110	0.3
c ₃	c	100	0.2
d ₁	d	10	1

Probabilistic DB



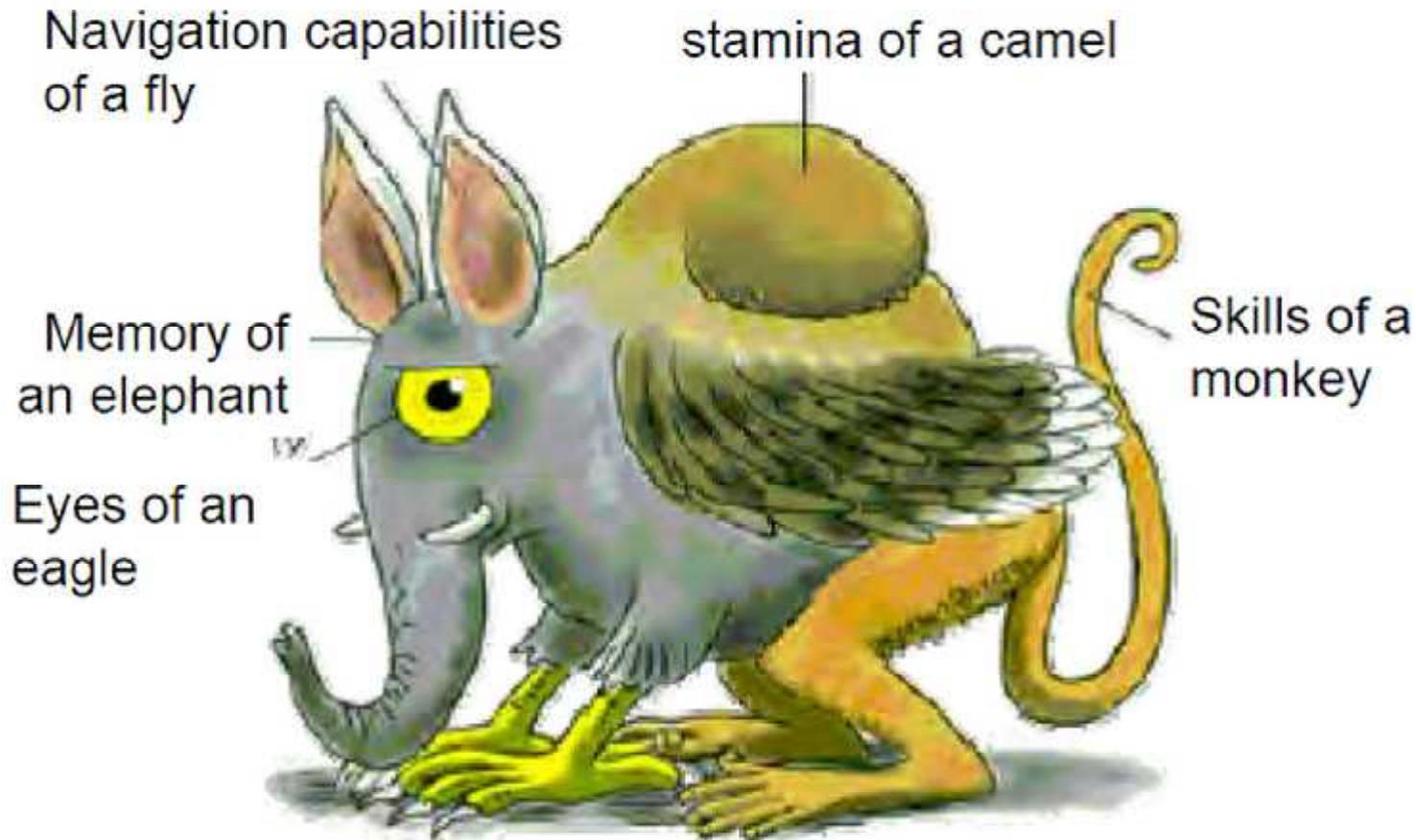
DB Community's Love Affair

SPEED, SPEED, SPEED
SCALE, SCALE, SCALE
FUNCTIONALITIES!!





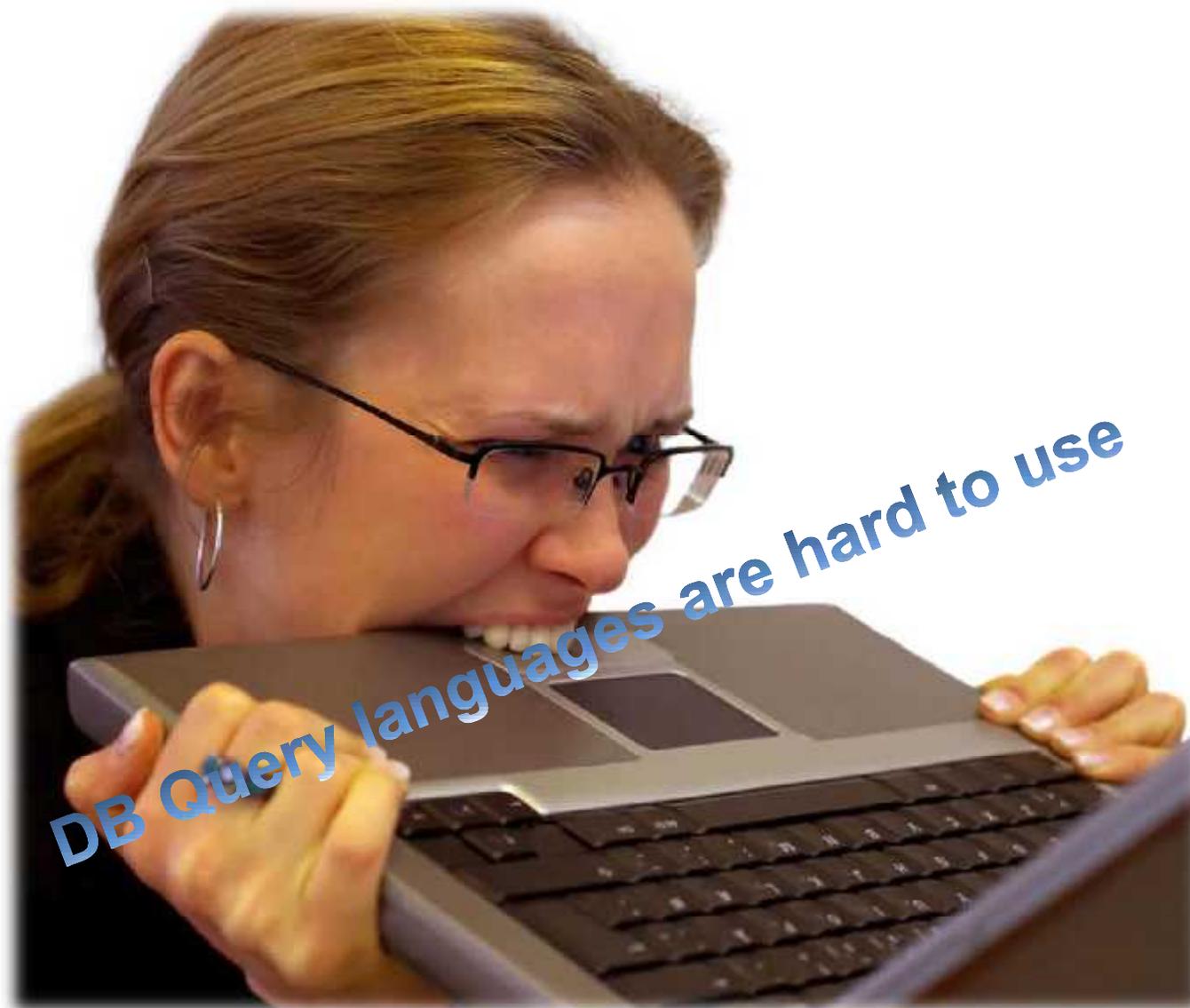
Assumptions Made By (Most) DM Systems



"The Perfect User"



The Real User





Query Formulation using SPARQL

```
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
PREFIX rss: <http://purl.org/rss/1.0/>
PREFIX dc: <http://purl.org/dc/elements/1.1/>

SELECT ?title ?known_name ?link
FROM http://planetrdf.com/index.rdf
FROM NAMED <phil-foaf.rdf>
WHERE {
  GRAPH <phil-foaf.rdf> {
    ?me foaf:name "Phil McCarthy".
    ?me foaf:knows ?known_person .
    ?known_person foaf:name ?known_name .
  }.
  ?item dc:creator ?known_name .
  ?item rss:title ?title .
  ?item rss:link ?link .
  ?item dc:date ?date.
}
ORDER BY DESC[?date] LIMIT 10
```





Reality Check!

Reality

“ Thirty years of research on query languages can be summarized by: we have moved from SQL to XQuery. At best we have moved from one declarative language to a second declarative language with roughly the same level of expressiveness. **It has been well documented that end users will not learn SQL; rather SQL is notation for professional programmers.**”

The Lowell Database Research Self-Assessment,
Communication of the ACM (May 2005)



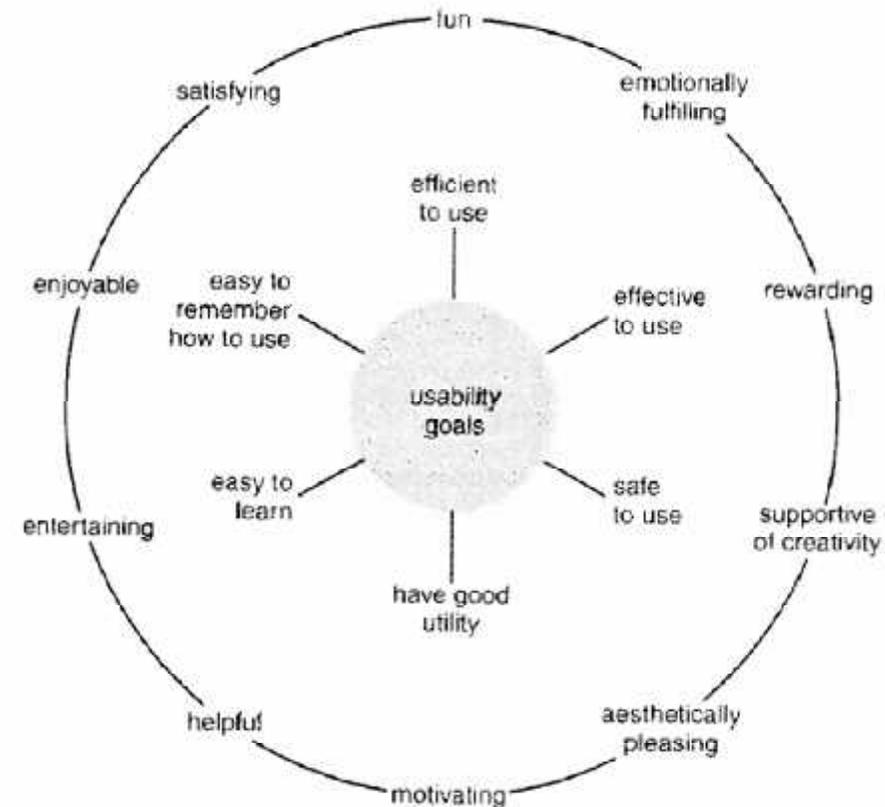
Usability [Preece et al.]

What is it?

How well users can use the system's functionality

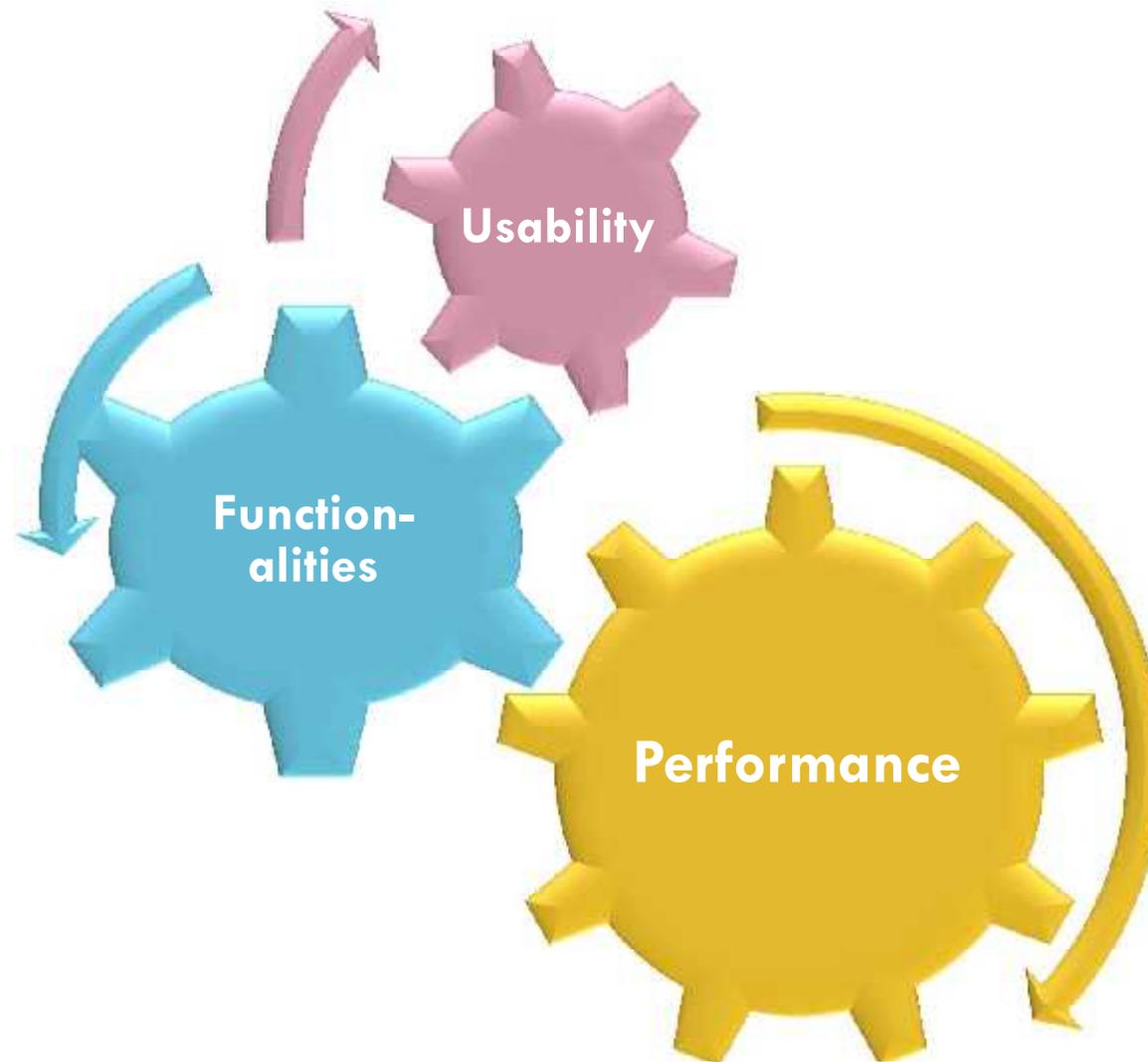
Dimensionality

- **Learnability:** is it easy to learn?
- **Efficiency:** once learned, is it fast to use?
- **Memorability:** is it easy to remember what you learned?
- **Errors:** are errors few and recoverable?
- **Satisfaction:** is it enjoyable to use?





DB Research Since 2006





Next..



- The World of HCI



What is Human Computer Interaction (HCI)?



“concerned with the design, evaluation and implementation of interactive computing systems for human use and with the study of major phenomena surrounding them”

of major phenomena surrounding them,
systems for human use and with the study

ACM SIGCHI (1992)



Usability and HCI

Studying and improving usability is part of Human-Computer Interaction (HCI).

Usability and good UI design are closely related





“The interface is the system.”

View

- Interface provides/conveys the **only** view of the “underlying” system

Perception

- Usable software sells better

Superficiality

- Users blame themselves for UI failings
- People who make buying decisions are not always end-users



Should and Must Do It Right

Ben Shneiderman

“Always should have “good” interfaces.
Computing time (power) is getting
cheaper but users’ time isn’t..”

“If the user can’t use it, it doesn’t work”

Susan Dray,
Distinguished Engineer of ACM





The World of Cool Interfaces!



Role of visual interfaces

“A picture is worth a thousand words. An interface is worth a thousand pictures

Ben Shneiderman, 2003

Google



Easy-to-use “dummy” visual interfaces are key to the spread of data management tool to wider community.

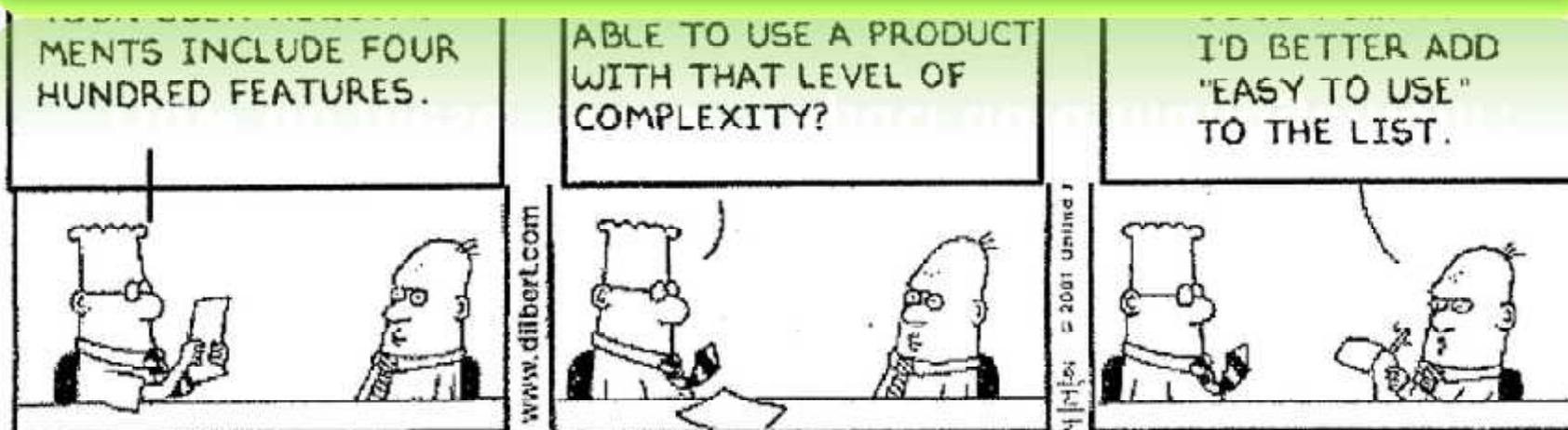




Lessons from HCI: Schneiderman's 8 Golden Rules

- Strive for consistency.
- Give shortcuts to the user.
- Offer informative feedback.
- Make each interaction with the user yield a result.
- Offer simple error handling.
- Permit easy undo of actions.
- Let the user be in control.
- Reduce short-term memory load on the user.

How do these “rules” impact data management?





Summary

“A USER INTERFACE IS LIKE A JOKE. IF YOU HAVE TO EXPLAIN IT, IT’S NOT THAT GOOD”

Anonymous, LinkedIn





The Chasm





Our Affinity to Visual Languages



Cave painting



Cuneiform



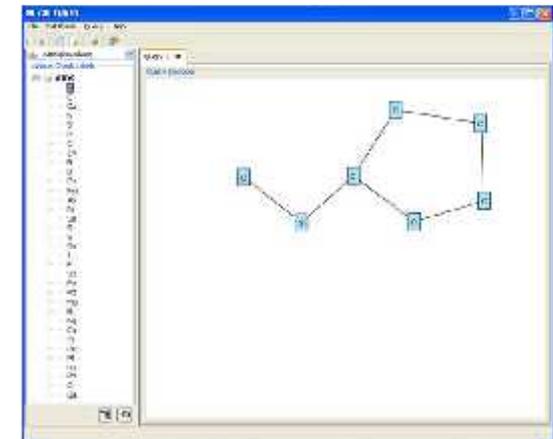
Hieroglyphics



Coats of arms



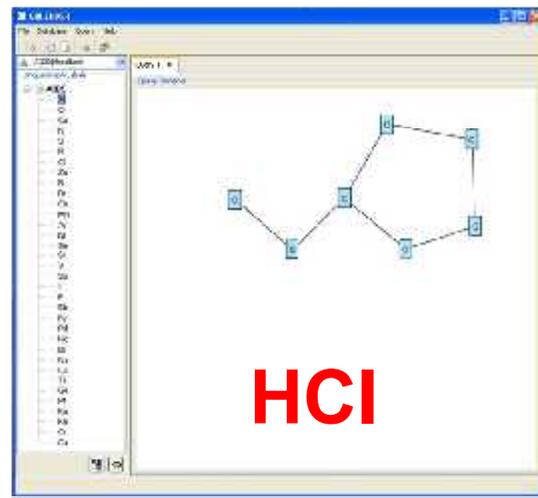
Emoji



Visual query language



Separate Ways for 40 years





Data Management Research





Querying Graph Databases

Query Formulation

- Formal query language
- GraphQL, SPARQL, PQL

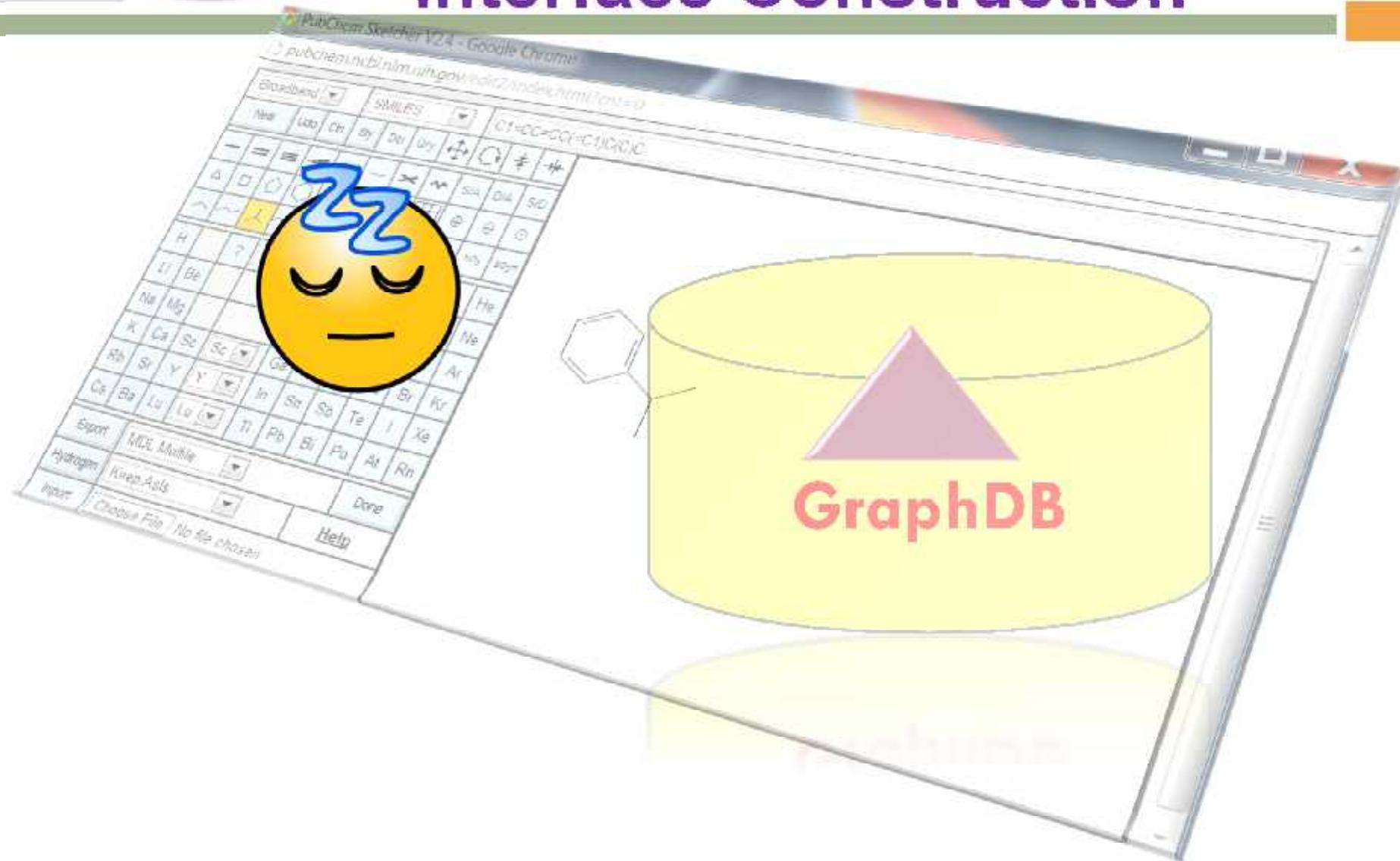


Query Processing

- Efficient algorithms and optimization techniques to process queries “quickly”
- FG-Index, Tree-PI, SSI, C-Tree, TALE



Classical Visual Query Interface Construction



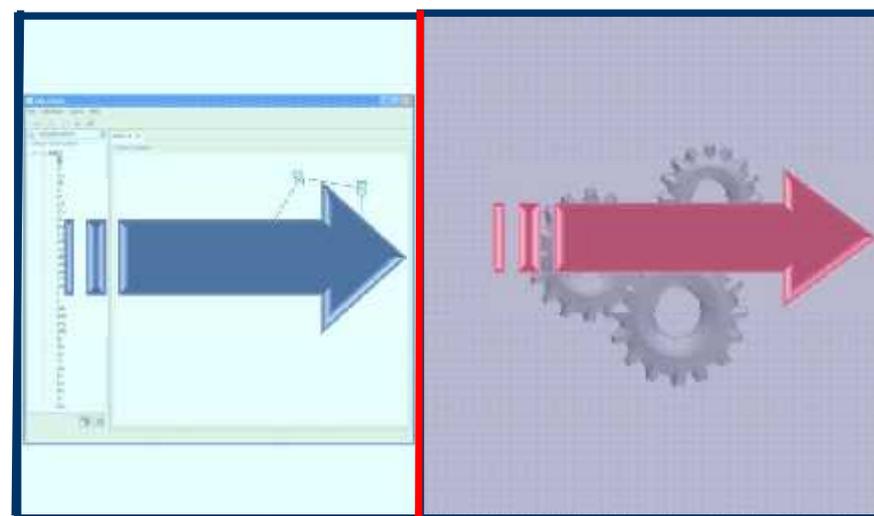
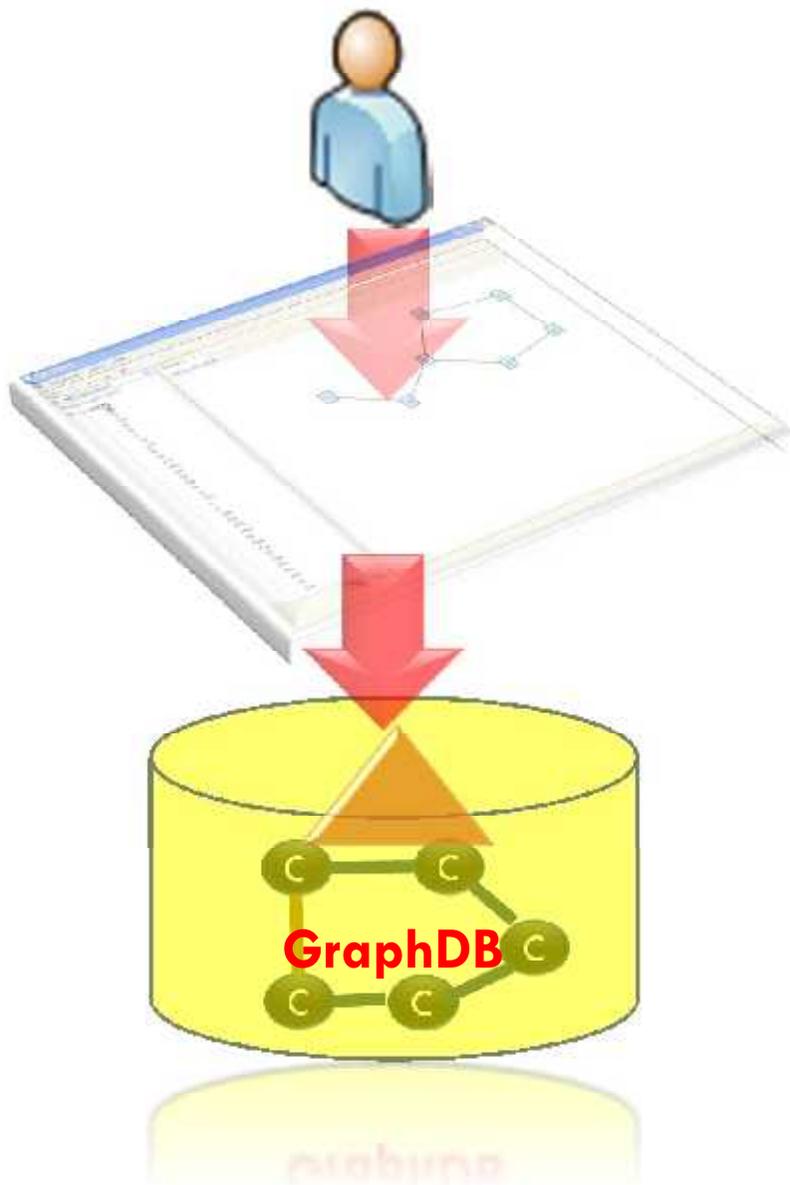


Classical Visual Query Interface Construction





Classical Visual Querying Paradigm



Query formulation Query processing

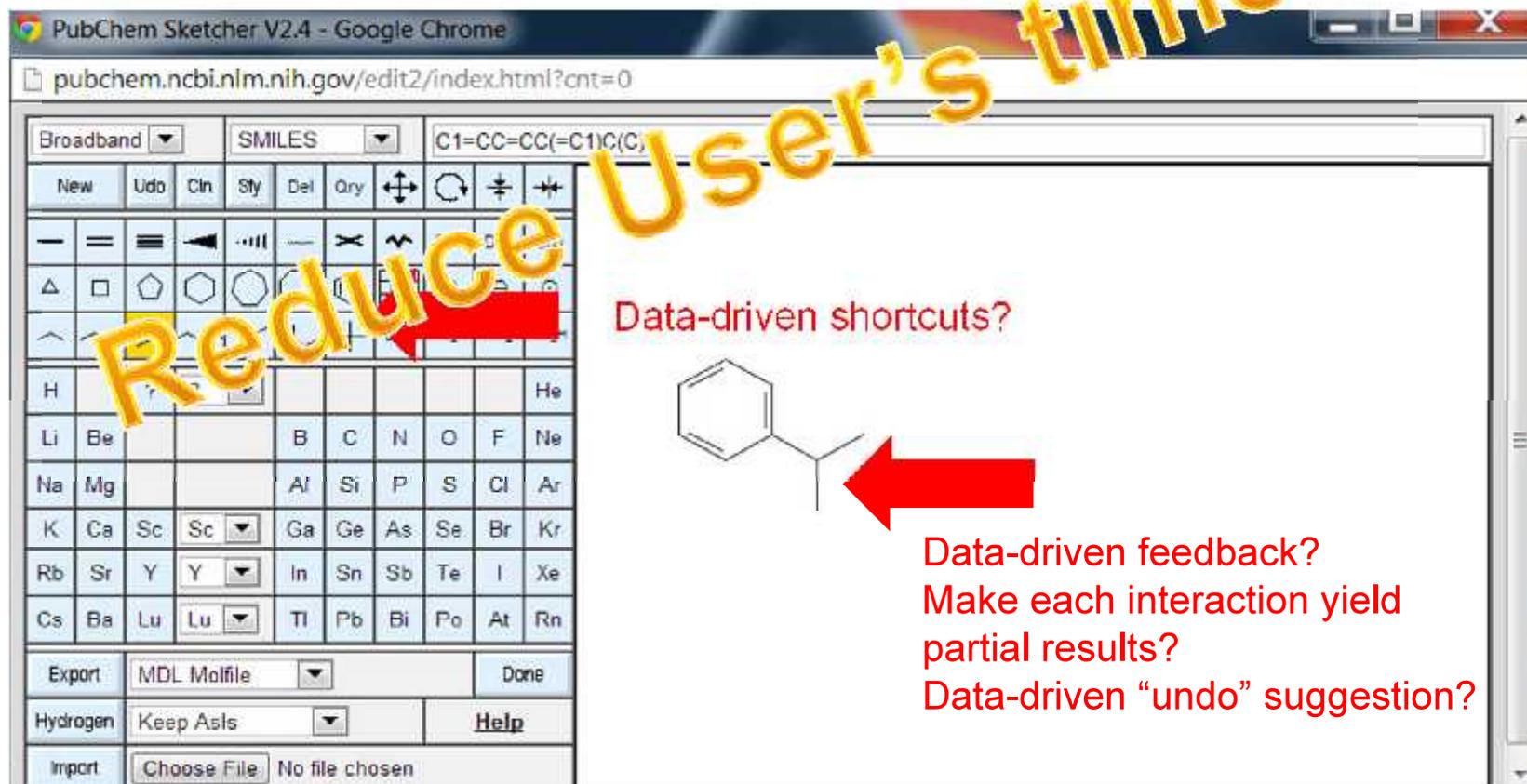
time →





What Happens if we Bridge the Chasm?

- Strive for consistency.
- Give shortcuts to the user.
- Offer informative feedback.
- Make each interaction with the user yield a result.
- Offer simple error handling.
- Permit easy undo of actions.
- Let the user be in control.
- Reduce short-term memory load on the user.



PubChem Sketcher V2.4 - Google Chrome

pubchem.ncbi.nlm.nih.gov/edit2/index.html?cnt=0

Broadband SMILES C1=CC=CC(=C1)C(C)

New Undo Cut Style Delete Copy Paste

Reduce User's time?

Data-driven shortcuts?

Data-driven feedback?
Make each interaction yield partial results?
Data-driven "undo" suggestion?

Export MDL Molfile Done

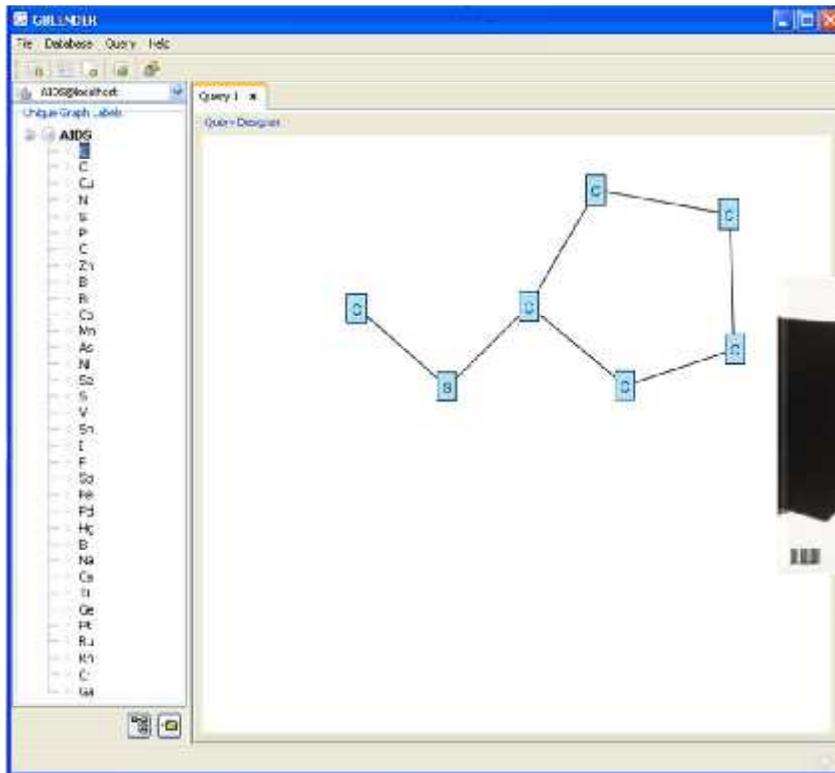
Hydrogen Keep AsIs Help

Import Choose File No file chosen

The screenshot shows the PubChem Sketcher interface. A large yellow diagonal text 'Reduce User's time?' is overlaid. A red arrow points from the text 'Reduce' to the 'Copy' button in the toolbar. Another red arrow points from the text 'Data-driven shortcuts?' to a chemical structure of 1-phenylethane. A third red arrow points from the text 'Data-driven feedback?' to the same chemical structure. The interface includes a SMILES input field with the string 'C1=CC=CC(=C1)C(C)', a periodic table, and various toolbars for editing and exporting.



HCI-aware Data Management



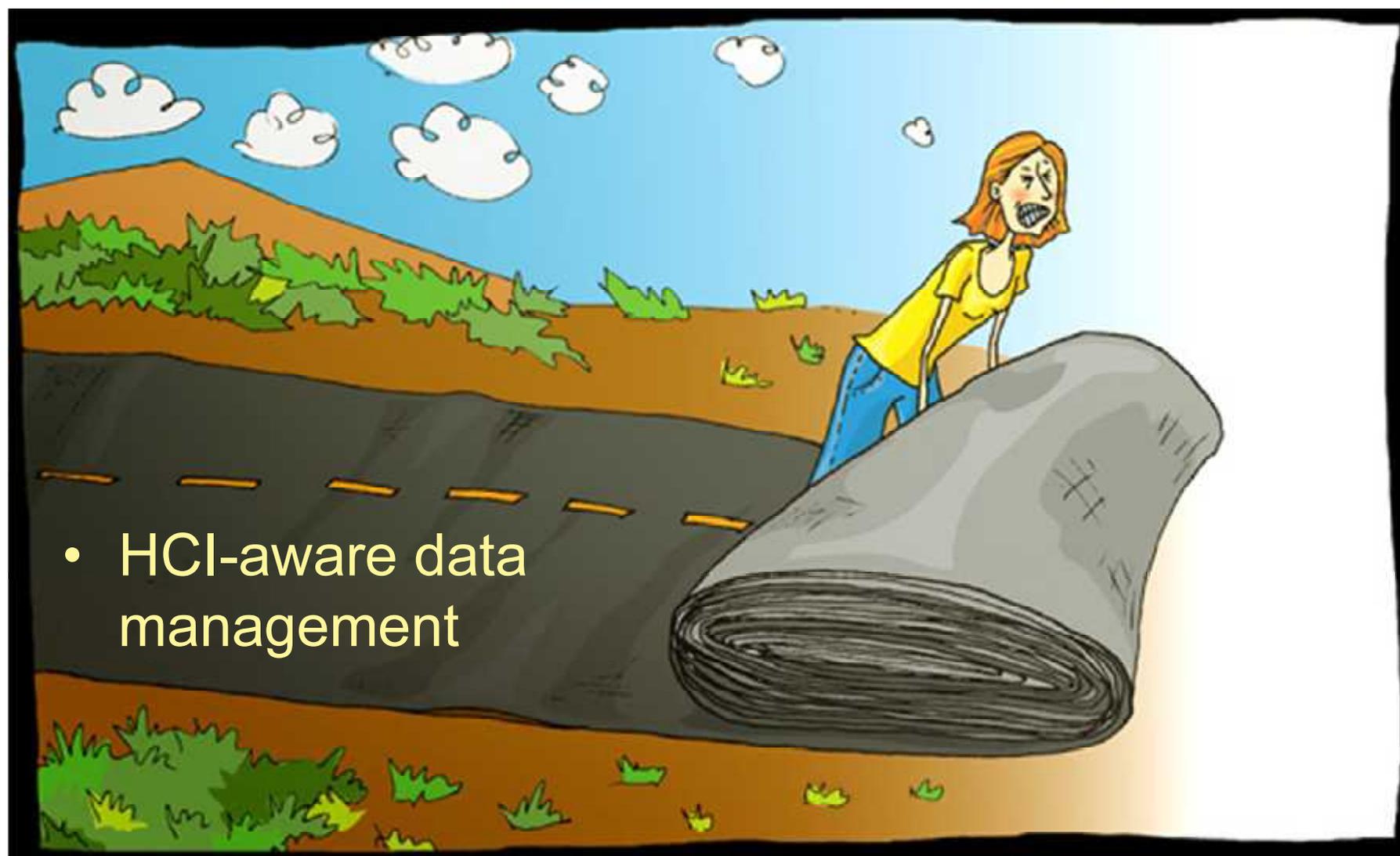
HCI



DB



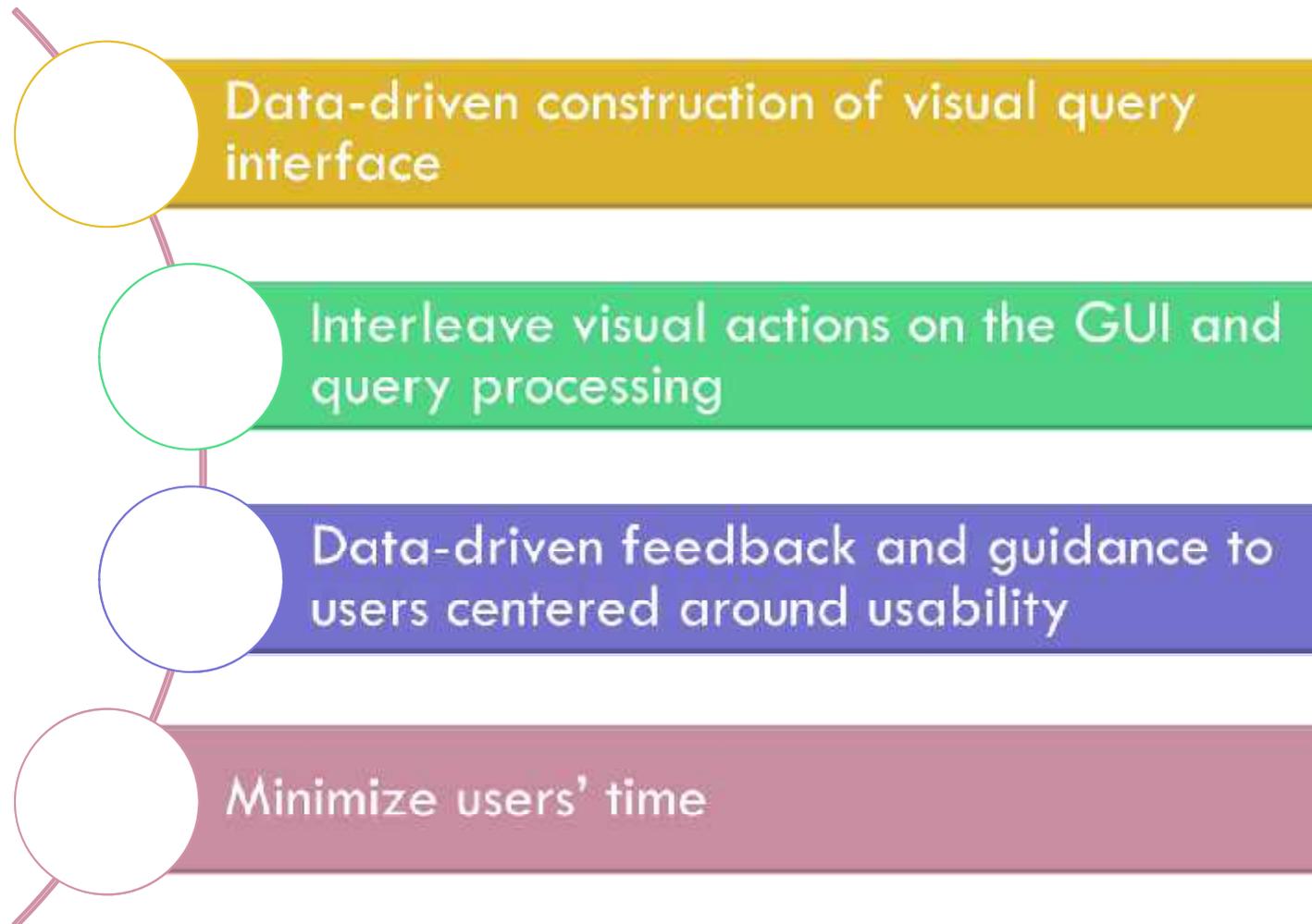
Next..



- HCI-aware data management

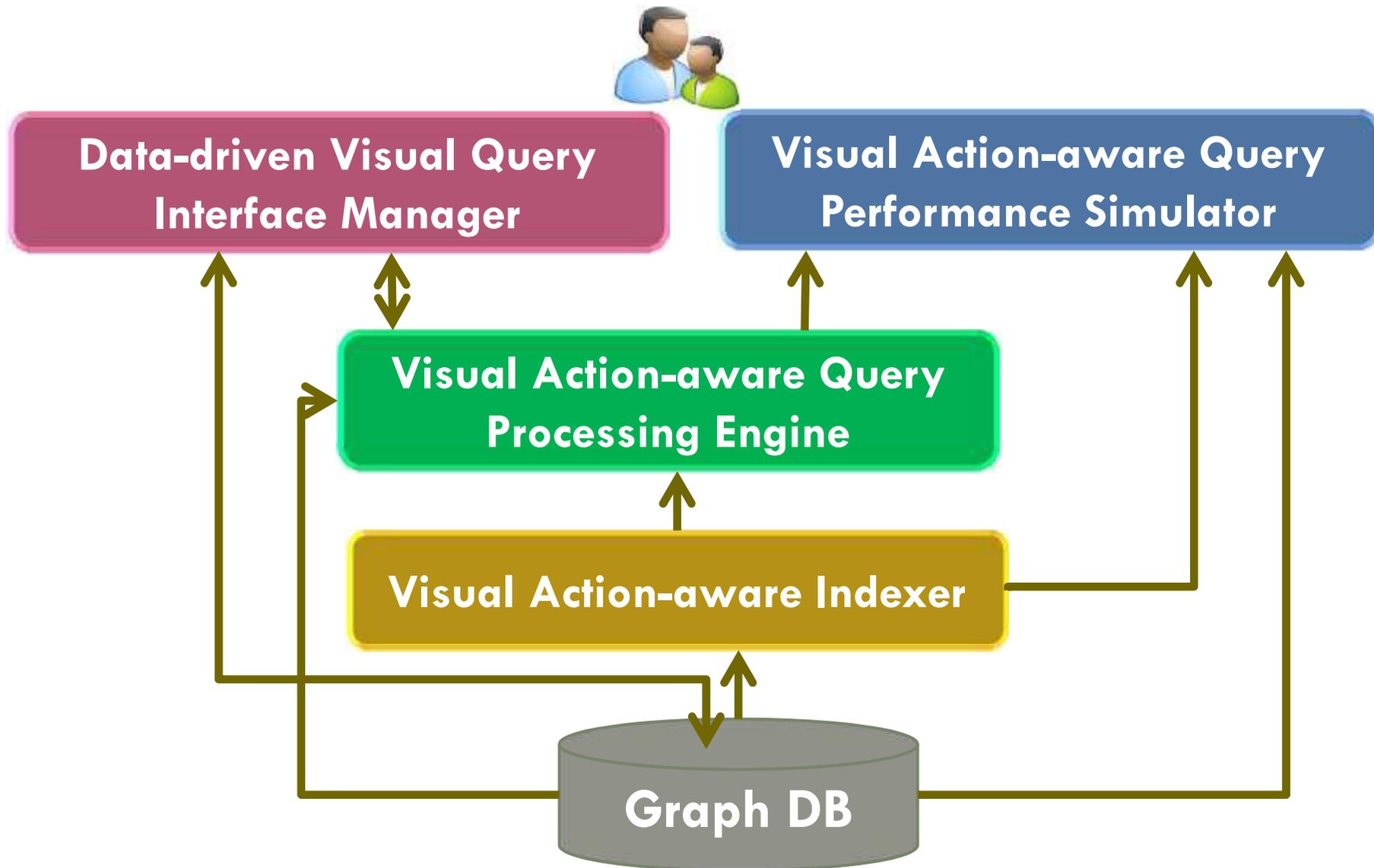


Broad Goals





Functional Architecture

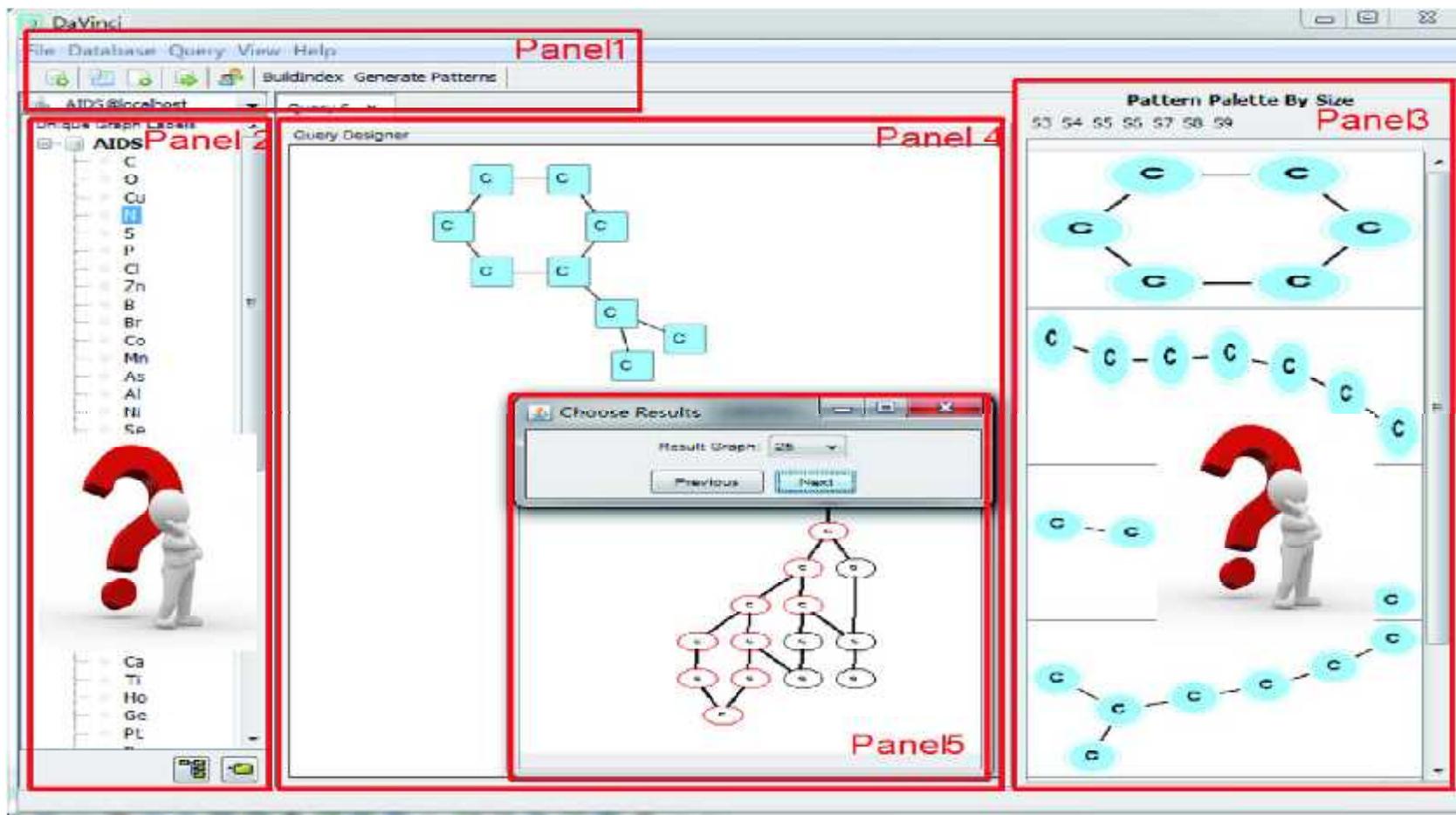




Data-driven Visual Interface Management

Goal 1

- Data-driven construction of visual query interface



The screenshot displays the DaVinci software interface, which is used for data-driven construction of a visual query interface. The interface is divided into several panels:

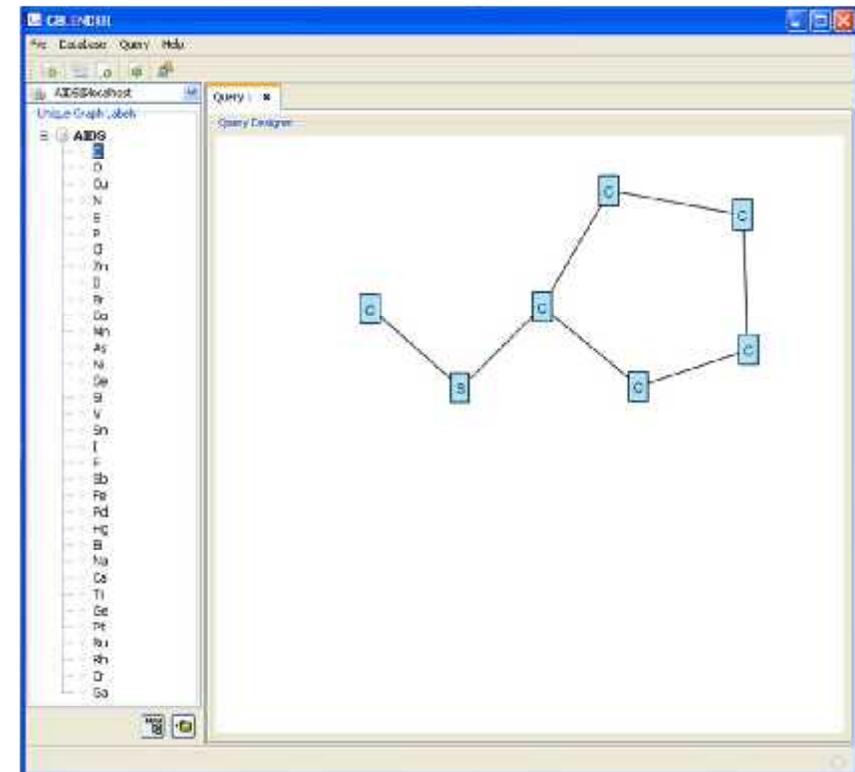
- Panel 1:** The top menu bar, including "File", "Database", "Query", "View", and "Help".
- Panel 2:** A list of elements (C, O, Cu, S, P, Cl, Zn, B, Br, Co, Mn, As, Al, Ni, Sp) and a 3D model of a person standing next to a large red question mark.
- Panel 3:** A "Pattern Palette By Size" window showing various molecular patterns (graphs) of carbon atoms (C) and a 3D model of a person standing next to a large red question mark.
- Panel 4:** The "Query Designer" window, which displays a graph of carbon atoms (C) and a "Choose Results" dialog box. The dialog box shows a "Result Graph" of 25 and "Previous" and "Next" buttons.
- Panel 5:** A window showing a hierarchical tree structure of carbon atoms (C).



Challenges

Panel 2

- Alphabetical, semantic, unordered
- How do we place the items for querying?
 - Search time to the same target position increases as the number of items in the panel increases [Byrne et al., CHI 99]
 - Semantic and alphabetical is faster [Halverson et al., CHI 08]
 - Users are faster at selecting targets that are closer to the top [Cockburn et al, CHI 07]
 - Last item effect [Bailly et al, CHI 14]

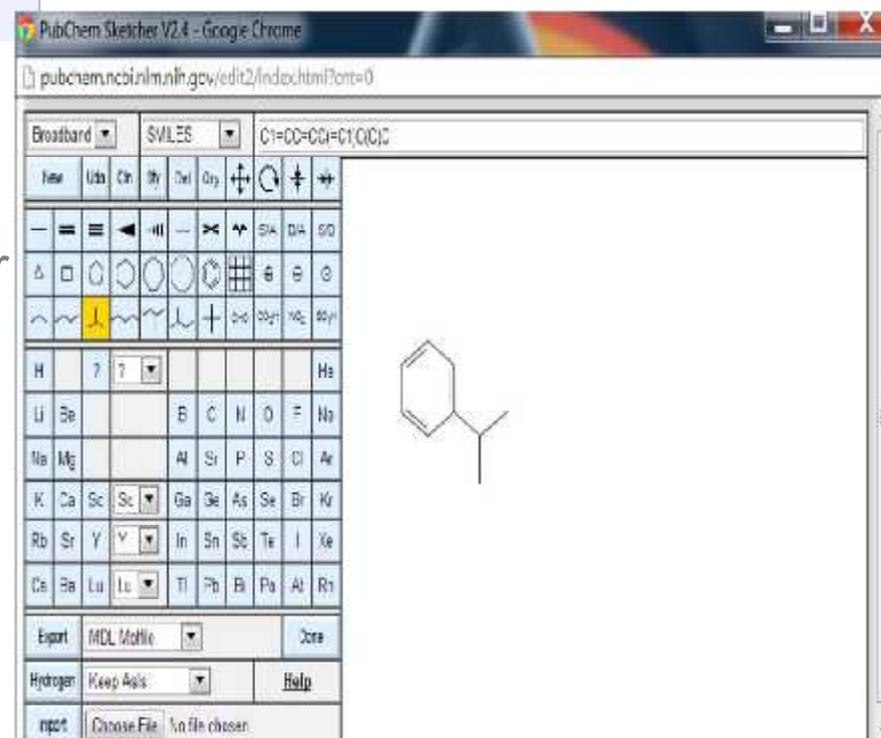




Challenges

Palette Selection

- Which subgraph patterns should be in the palette?
 - Formulate query easily and faster
 - Give shortcuts
- Issues
 - “Data-driven” selection
 - Size of the palette
 - Frequent and infrequent patterns
 - Maximally covers the DB
 - Minimal redundancy among patterns
 - Dynamically maintained

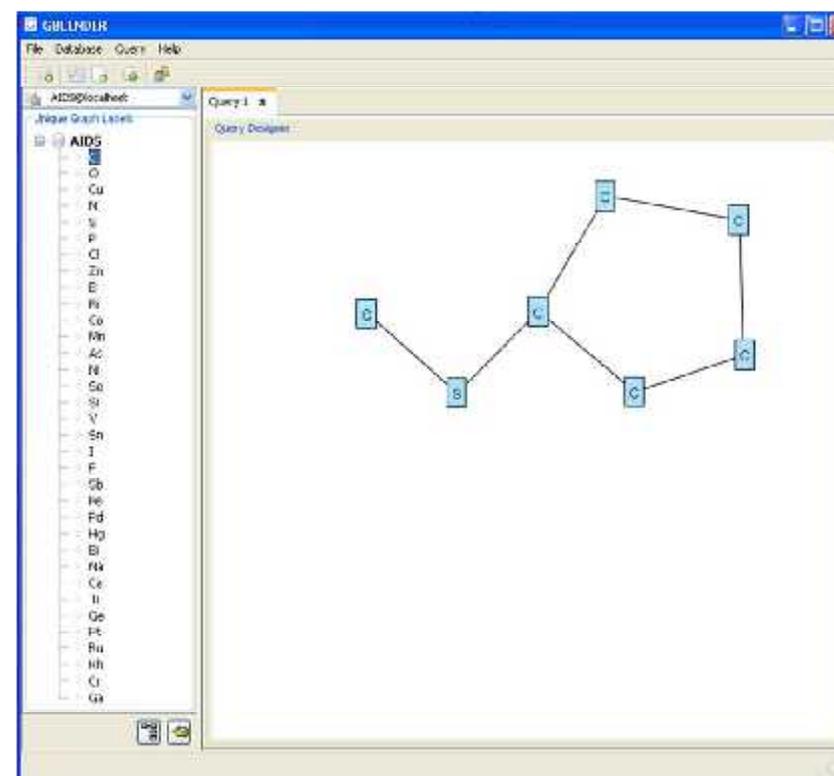




Visual Query Formulation

Characteristics

- **Node/edge-at-a-time** or **subgraph-at-a-time** approach
 - Size of the query fragment increase by **1** or **k**
- Query can be modified at any time
 - Size may not increase monotonically
- Different formulation sequences
- At any step, the partial query is either **frequent** or **infrequent**
 - The chance of a fragment to remain frequent diminishes as the size of the query grows

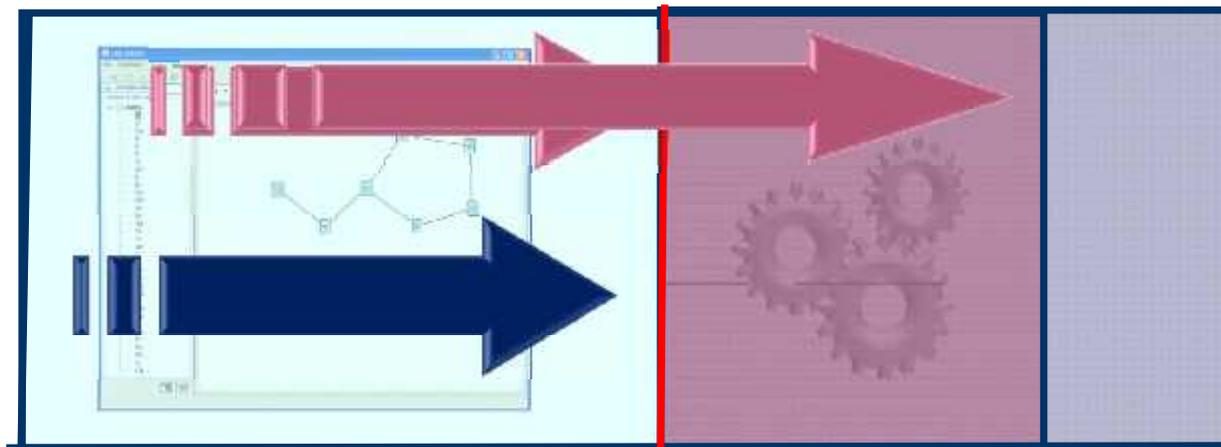




Visual Action-aware Query Processing

A novel paradigm

- Why **wait** for the complete visual query to be constructed **before** initiating query evaluation? How can we blend these two steps?
- By initiating query processing “early”, can we significantly **reduce** the **user’s waiting time**?



Query formulation

Query processing

time





Non-traditional Challenges

Action-aware Indexing

- Prune irrelevant results even when **partial** query graph is known
- Efficient traversal from g' (Step i) to g'' (Step $i+1$) where $g' \subset g''$ and $|g''| = |g'| + k$
- Support of pruning **graph-structured** frequent and infrequent fragments
- **Smaller-sized** graphs should be efficiently indexed
- **Minimize** candidate verification of partial results

Intermediate results materialization

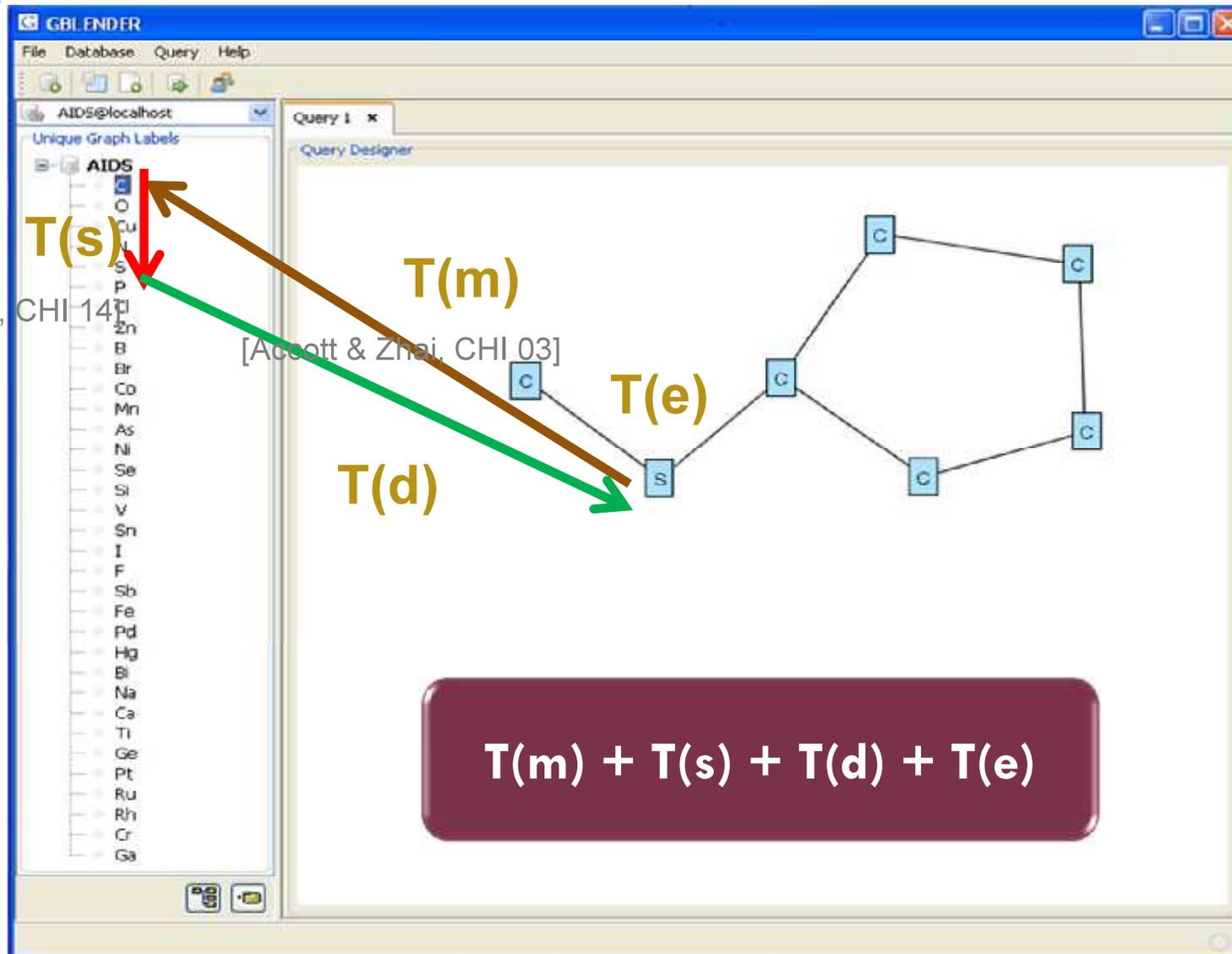
- Materialization of all partial candidate graphs matching the query fragment is necessary
- Unreasonable in traditional DB
- Challenging due to computational hardness of subgraph matching



Bounds on materialization time

[Bailly et al, CHI 14]

[Accott & Zhai, CHI 03]



$$T(m) + T(s) + T(d) + T(e)$$



Non-traditional Challenges

Selectivity-free query processing

- Selectivity-based query processing is impractical
- Query fragments are drawn in any arbitrary sequence
- Complete query is not available during query evaluation due to the paradigm
 - “Push-down” highly selective fragment is not possible
 - Classical approach of physical query plan generation is ineffective

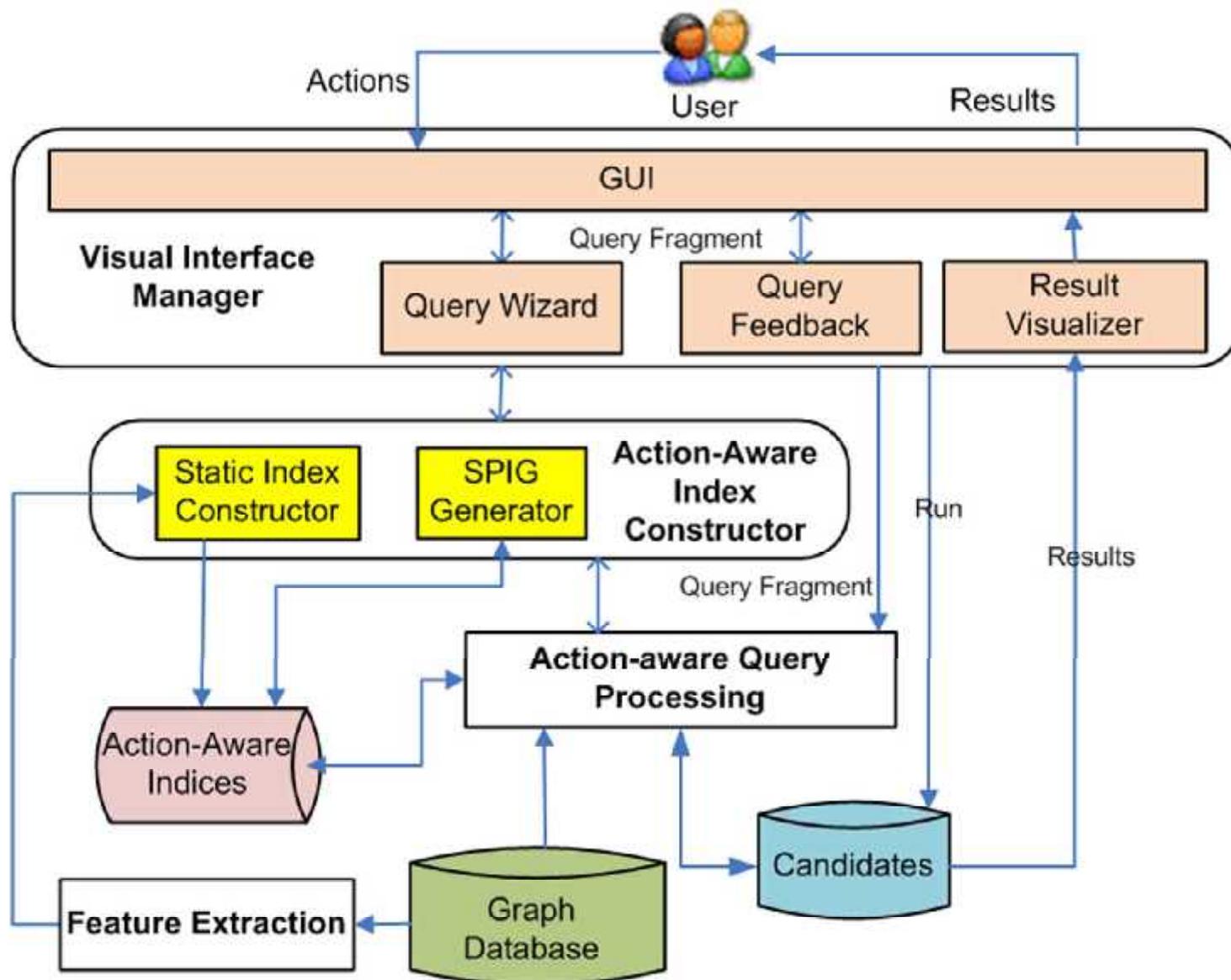
Focus on waiting time of users

- Quick formulation of queries
- **System response time (SRT)** matters more than backend processing cost!
- SRT should be **robust** to different query formulation sequences.

“Computing time (power) is getting cheaper but users’ time isn’t..”

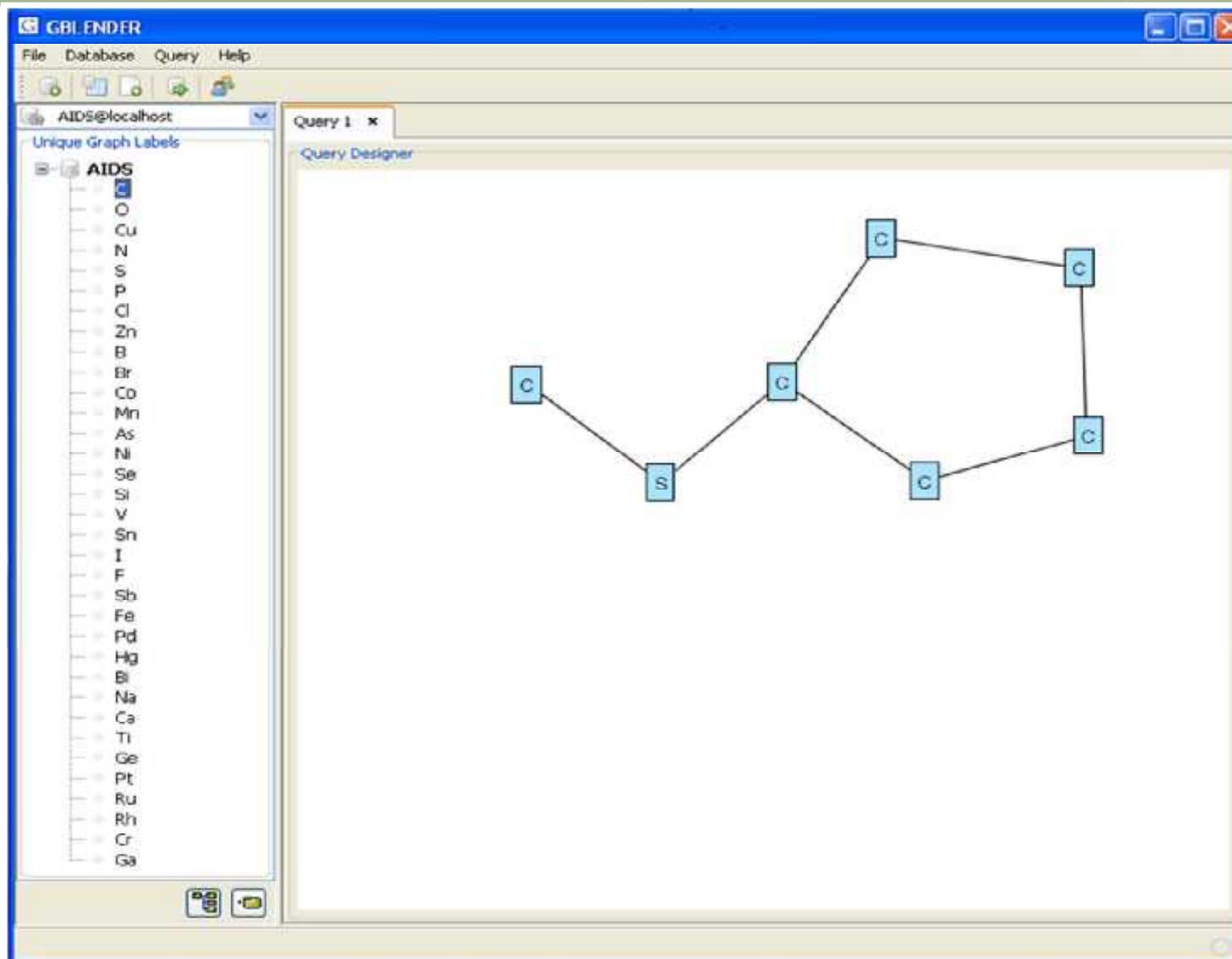


VOGUE [CIDR 13; SIGMOD 10, 11, ICDE 12]





A Video of Visual Action-aware Query Processing

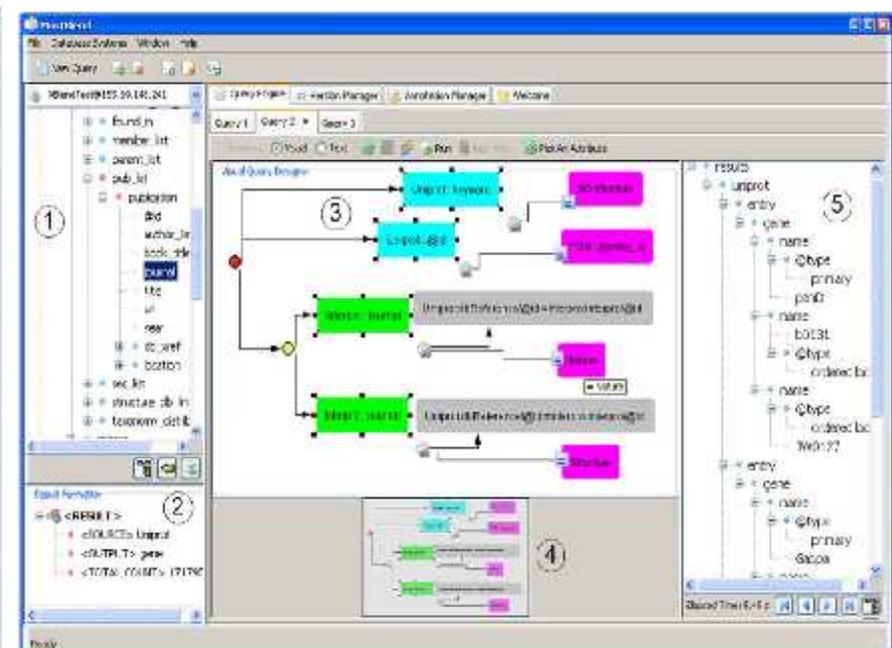
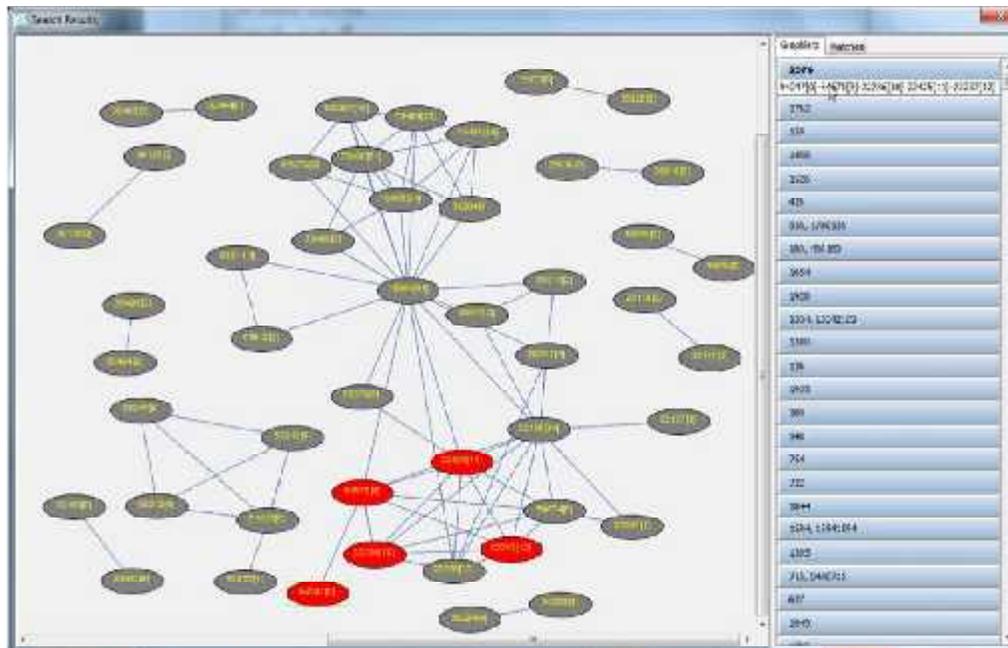




Related Efforts

- ☑ QUBLE [SIGMOD 13, VLDBJ 14]
 - @ For large networks

- ☑ MUSTBLEND [DASFAA 13, ICDE 09, ICDE 06]
 - @ For XML data





Data-driven Visual Interface Management Revisited

Goal 2

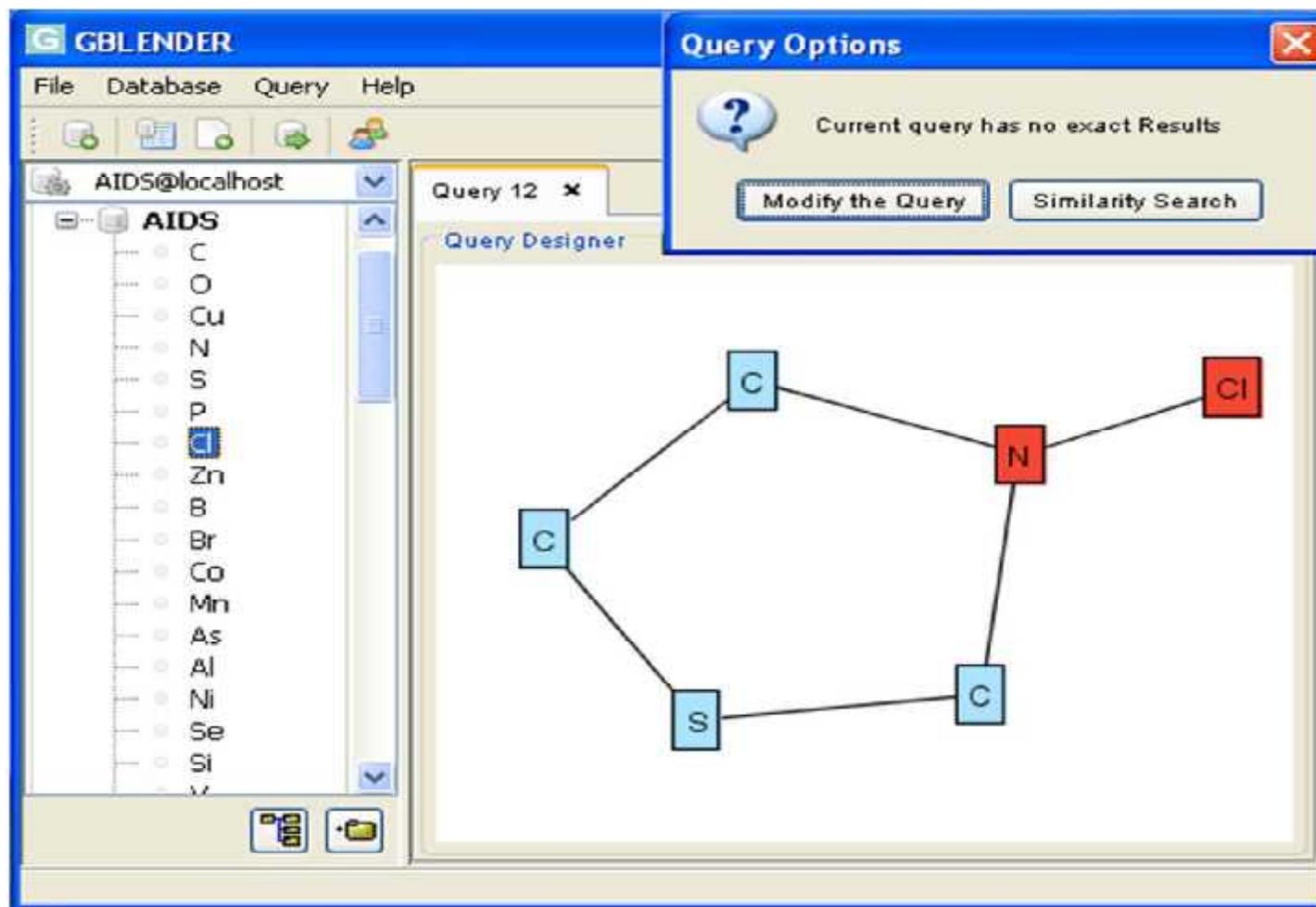
- ❑ Side-effects of intermediate results materialization
- ❑ Data-driven feedbacks and notifications

What kind of feedback?

- Aiding query construction by making appropriate suggestions
 - Given a query fragment, which top-k patterns are most likely to be added to the query?
- Empty results detection



Empty Results Detection



The screenshot displays the GBLENDER software interface. On the left, a list of elements is shown under the 'AIDS' database, with 'C' (Carbon) selected. The main window shows a 'Query Designer' with a molecular graph consisting of six nodes: five blue 'C' nodes and one red 'N' node. The nodes are connected as follows: a top 'C' node is connected to a middle-left 'C' node and a middle-right 'N' node; the middle-left 'C' node is connected to a bottom-left 'S' node; the bottom-left 'S' node is connected to a bottom-right 'C' node; the middle-right 'N' node is connected to a rightmost 'Cl' node and the bottom-right 'C' node. A 'Query Options' dialog box is open in the top right, displaying a question mark icon and the text 'Current query has no exact Results'. Below this text are two buttons: 'Modify the Query' and 'Similarity Search'.



Data-driven Visual Interface Management Revisited

Goal 3

- ❑ Interruption-sensitive notifications

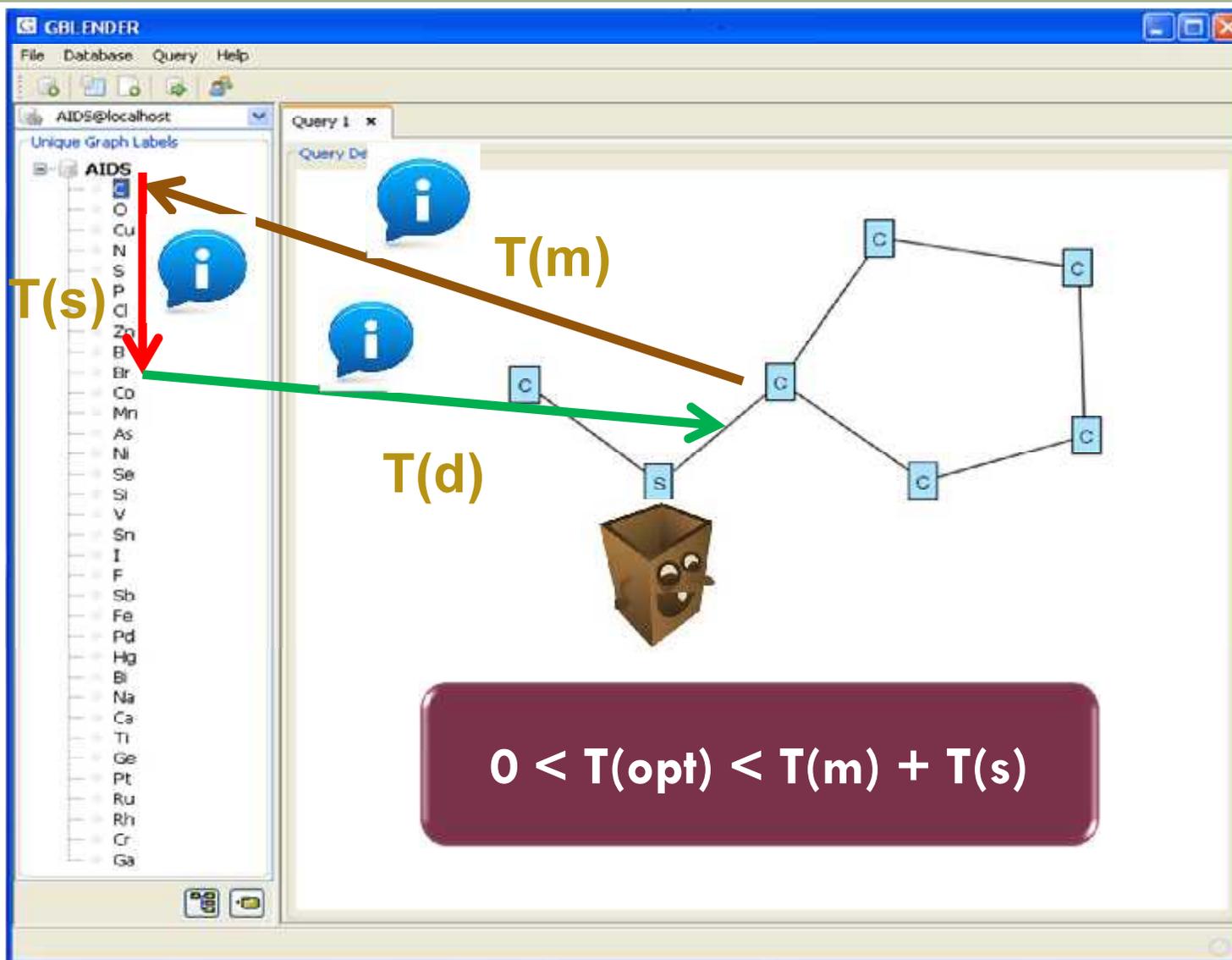
Lessons from HCI & Cognitive Psychology

- Interrupting users engaged in tasks by delivering notifications **inopportunistically** can negatively impact task completion time, lead to more errors, and increase user frustration [Bailey et al., JCSB 2006].
- Breakpoint-based notifications





Bounds on feedback generation time



The screenshot shows the GBLENDER application window. On the left, a list of 'Unique Graph Labels' includes elements like O, Cu, N, S, P, Cl, Zn, Br, Co, Mn, As, Ni, Se, Si, V, Sn, I, F, Sb, Fe, Pd, Hg, Bi, Na, Ca, Ti, Ge, Pt, Ru, Rh, Cr, and Ga. A red arrow labeled $T(s)$ points from the top of this list to a specific element. A green arrow labeled $T(d)$ points from that element to a node in a graph. A brown arrow labeled $T(m)$ points from the graph back to the top of the list. The graph consists of several nodes labeled 'c' and one node labeled 's' which is represented by a cartoon box character. A dark red rounded rectangle at the bottom contains the inequality:

$$0 < T(\text{opt}) < T(m) + T(s)$$



Visual Action-aware Query Performance Simulator

Large-scale performance study

- Traditional approach
 - Randomly extract subgraphs of different size and execute them
- Doesn't work in this paradigm!



Why?

- Queries need to be visually constructed by users
- GUI latency is critical for performance study

Challenge

- Users are expensive!
- How do we simulate visual query formulation?





Next..



- DB in the changing world
- HCI in the changing world
- The chasm!
- HCI-aware data management
- **Conclusions**



Blending more Complex Queries

Homeomorphic queries

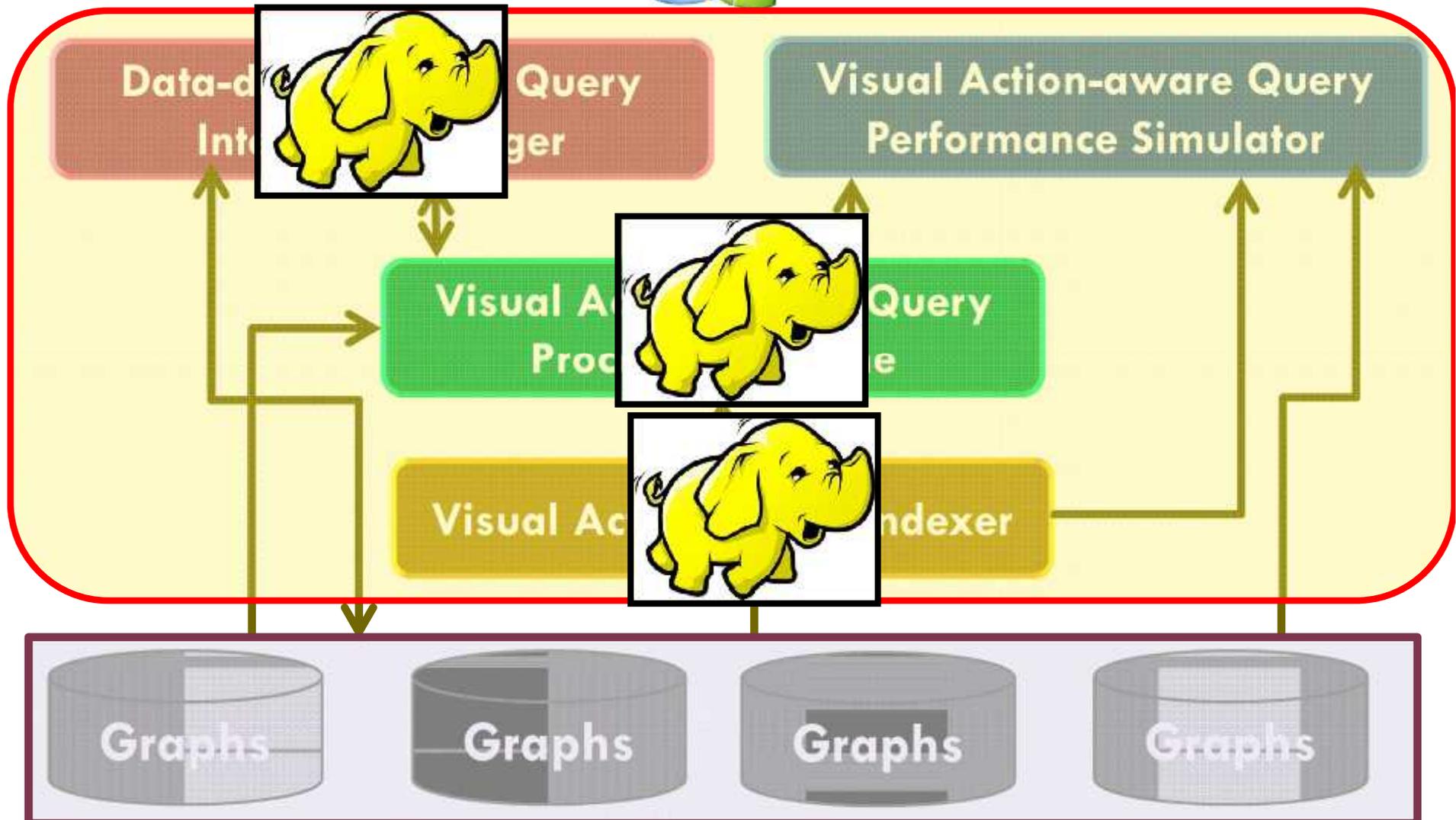
Distance-constraint queries

Multi-attribute queries

Graph simulation

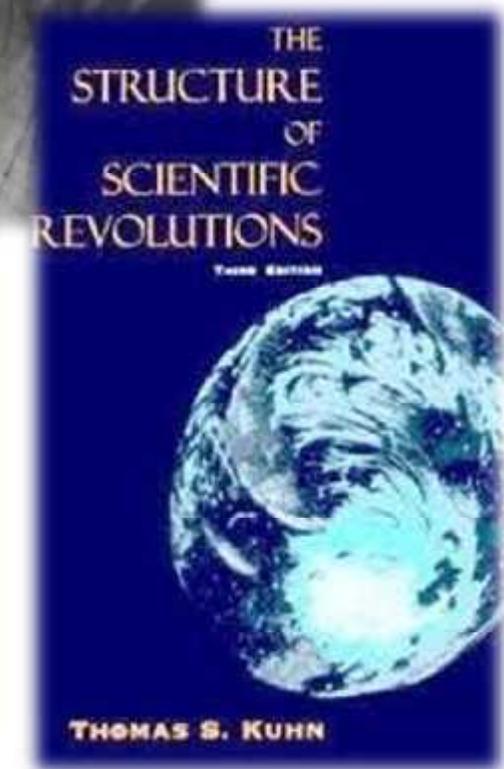


HCI-aware Data Management in the Big Data Era





Kuhn's Paradigm

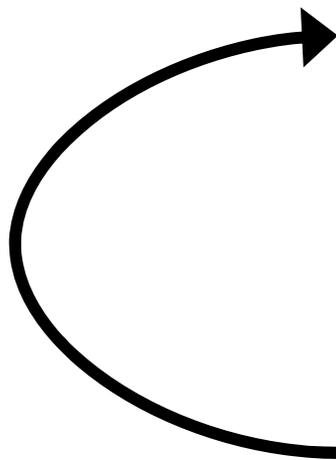


Published in 1962

1) Establish a paradigm

2) Mature science

3) Anomalies appear





Kuhnian Example



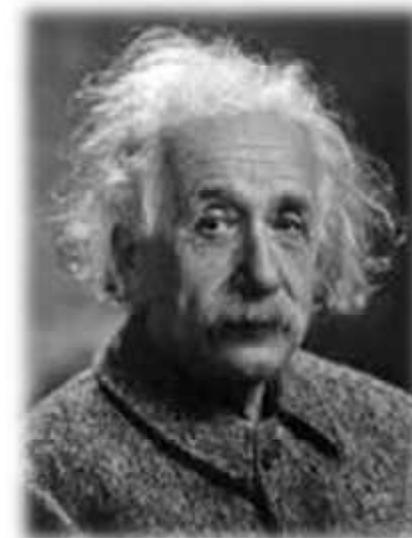
1700: **Classical (Newtonian) mechanics**

1800

Experiments with light show anomalous behavior

1900

1905: **Relativity**





Kuhnian Paradigm for Data Management



1970s: Query Formulation → Query Processing



2000s: Visual query formulation ↔ query processing

- **Blends** data management and HCI
- Three key components
 - Data-driven visual interface management
 - Visual action-aware query processing
 - Visual action-aware performance simulation
- Can be extended to wide variety of data types



Broad goal

- Pervasive desire to stimulate a cultural shift in our thinking by bringing together HCI and data management to “work” together



Acknowledgements



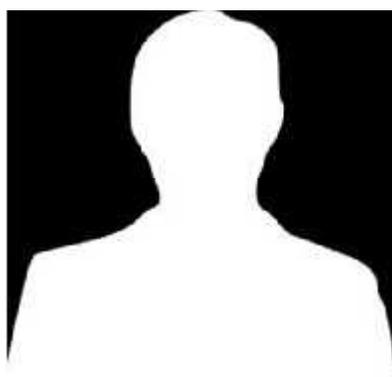
Byron Choi
Hong Kong Baptist Univ



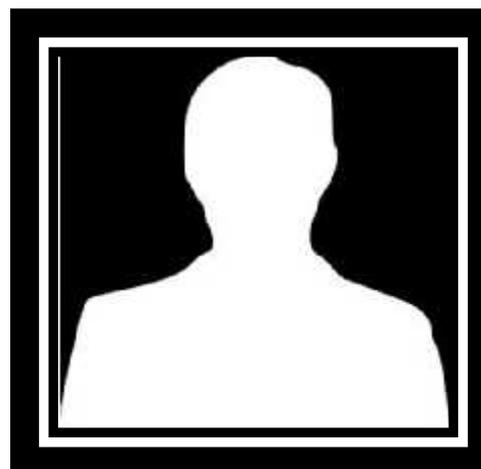
Changjiu Jin



Shuigeng Zhou
Fudan Univ, China



Ba Quan Truong



Ho Hoang Hung

