



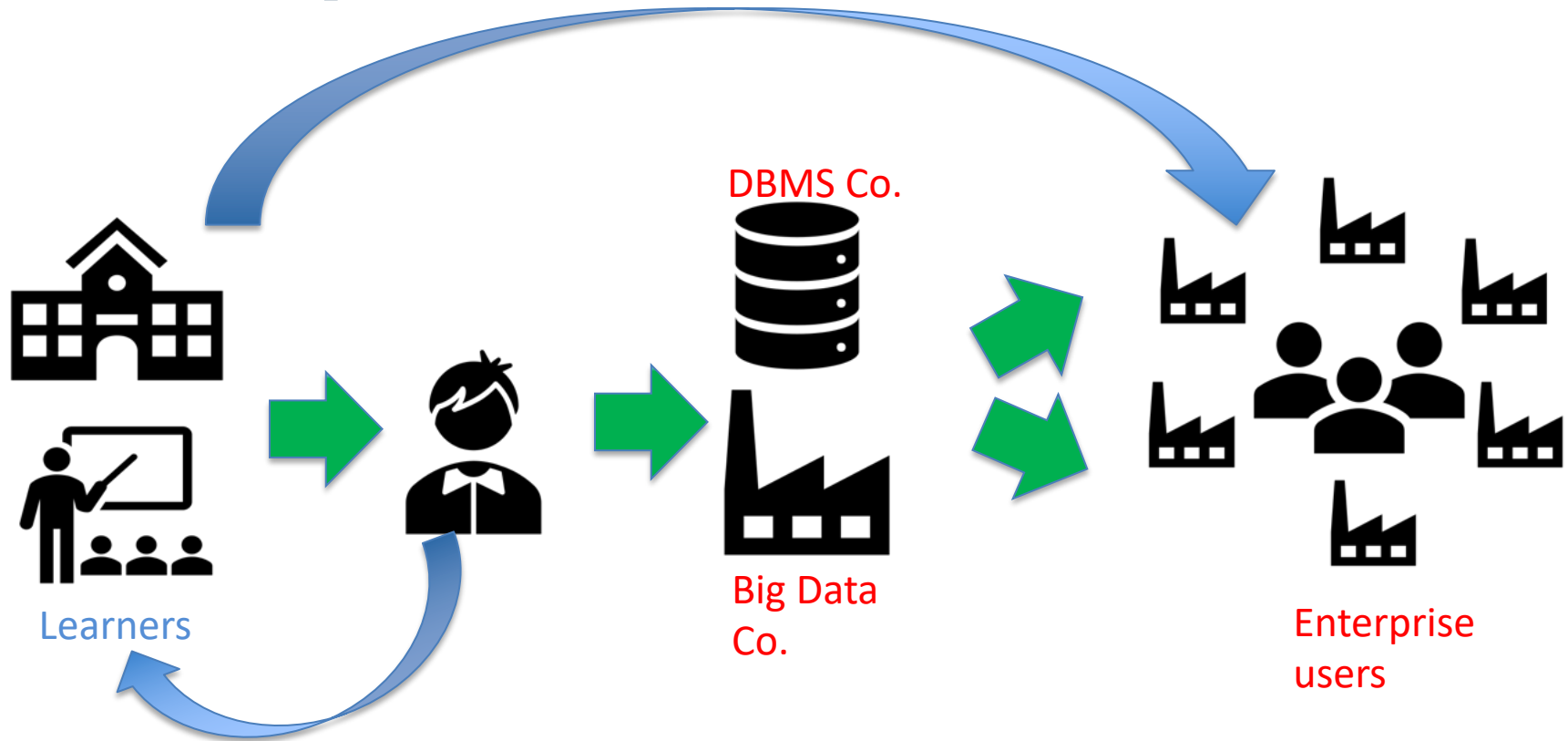
NANYANG
TECHNOLOGICAL
UNIVERSITY
SINGAPORE

Human Learners of Relational Query Processing: **Who Cares?**

Sourav S Bhowmick
assourav@ntu.edu.sg



DB for Enterprise Users & Developers



Beyond Enterprise Users?



students



biologist



Pharmacist/chemist



ecologist



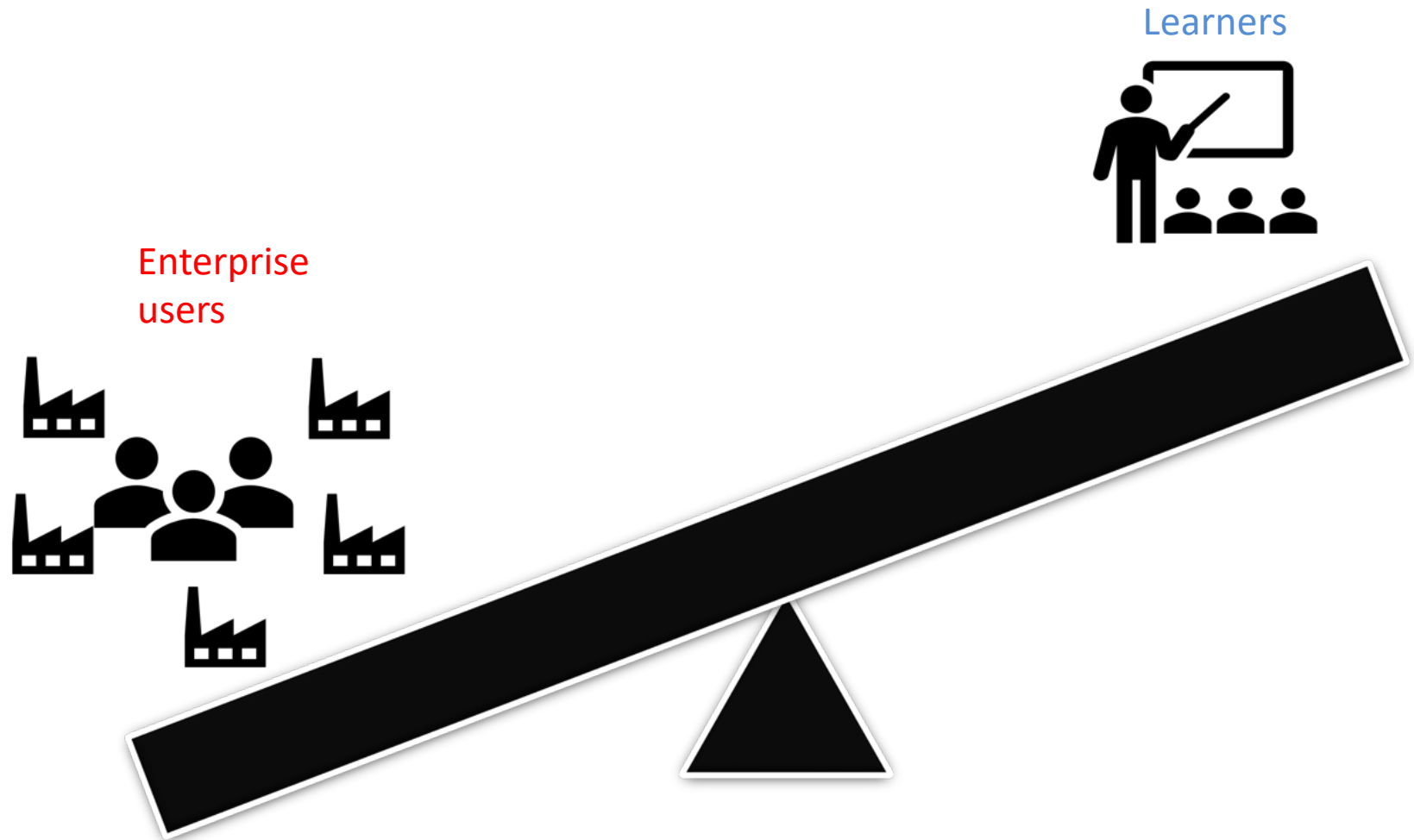
journalist



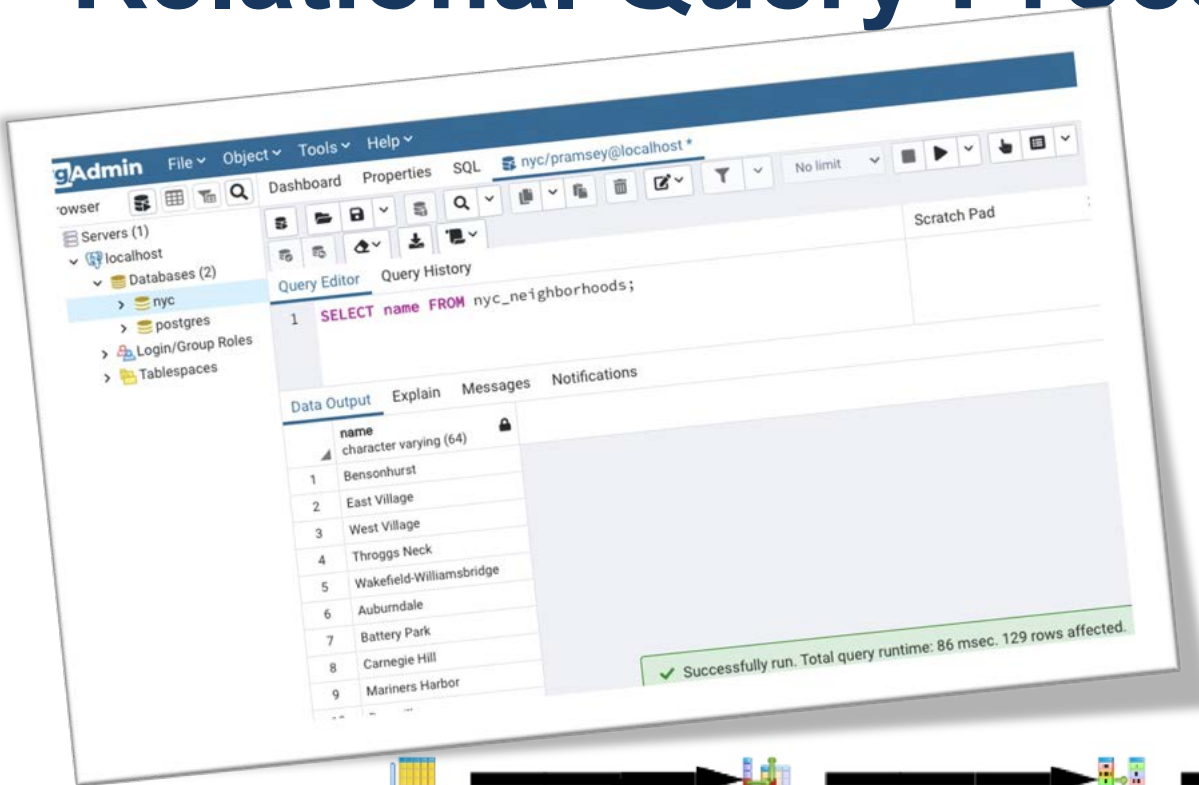
social scientist



Where Our Attention Lies?



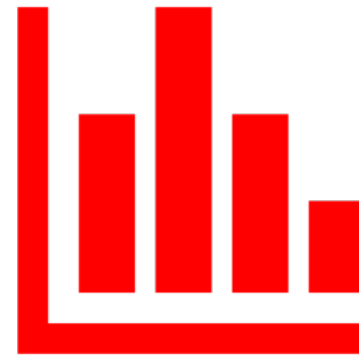
Scant Support For Learning Relational Query Processing



The Popularity of Data Science & AI



Data



Value



The Changing Landscape of Learning

How are you closing the skills gaps that exist in your organisation?

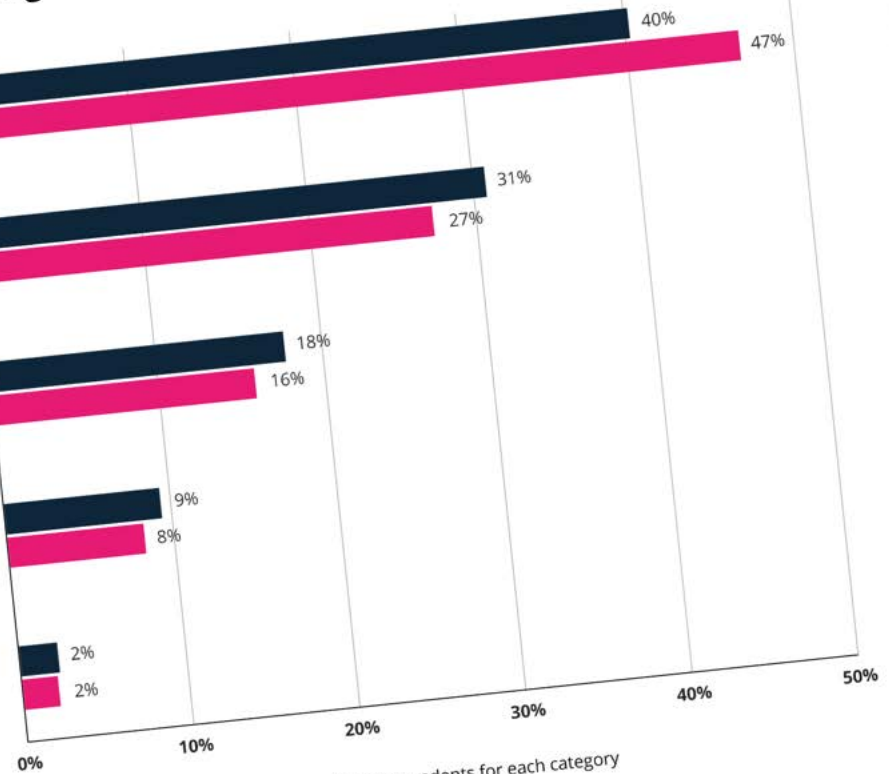
Upskilling employees
(updating existing
skills to meet
latest standards)

Reskilling employees
(training employees for
skills that traditionally
fell outside of their
scope of responsibility)

Hiring new employees
with specialised skills

Hiring part-time
contractors
(‘giggers’) to fulfil
specific skills gaps

Other



● HR, L&D, talent managers or a related role

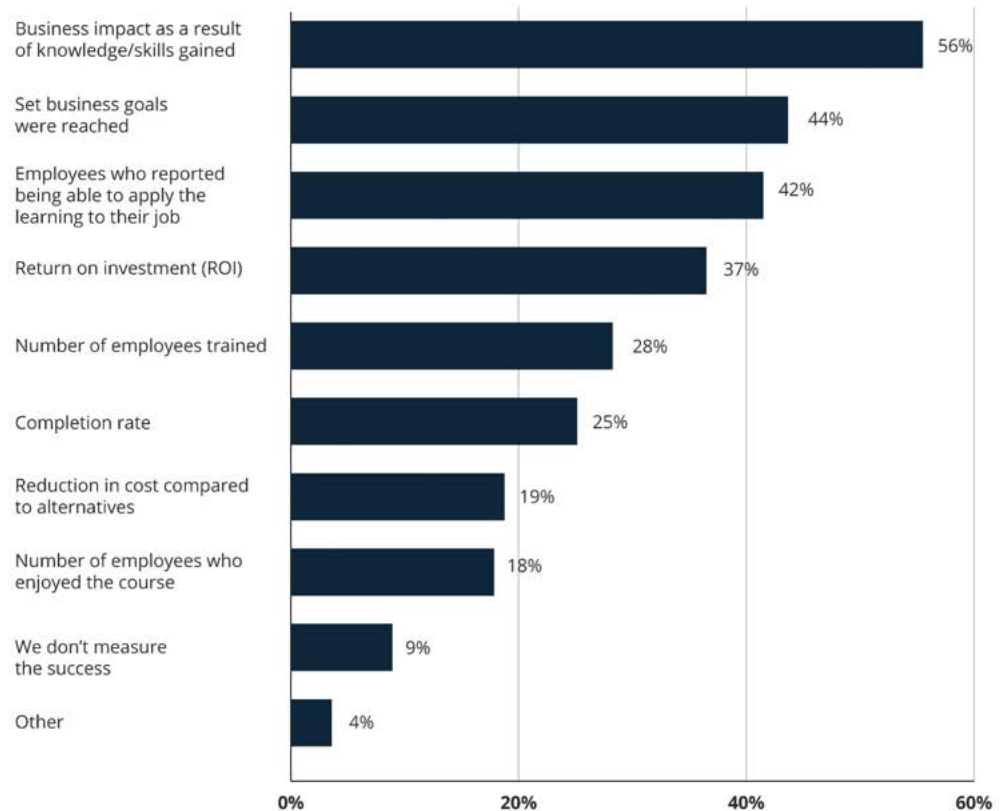
● People Managers

<https://www.getsmarter.com/blog/career-advice/why-its-important-for-corporates-to-encourage-lifelong-learning/>



The Changing Landscape of Learning

How does your organisation measure the success of its learning and development initiatives?



<https://www.getsmarter.com/blog/career-advice/why-its-important-for-corporates-to-encourage-lifelong-learning/>



Question 1

What are the observed challenges brought by the traditional modes of learning of relational query processing?



Core Topics of Relational Query Processing

Topics

- Set of physical operators
- Query processing models
- Selection of query execution plans
- Cost estimation of a query plan

Modes of Learning

- Seminars, lectures
- Textbooks, online resources
- Off-the-shelf RDBMS



Is Understanding Query Execution Plans (QEP) a Challenge?

Sem Y:

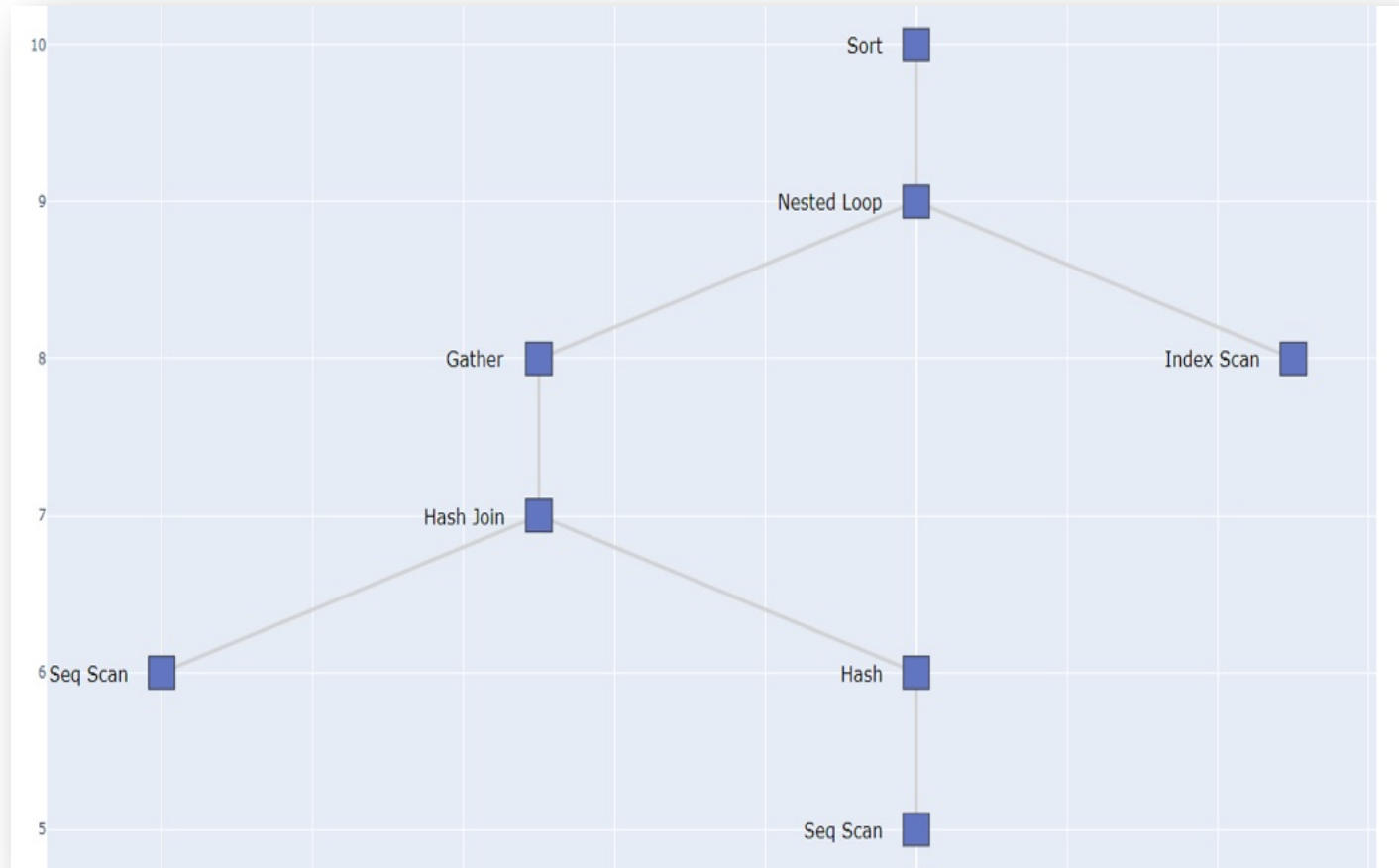
Class size: 162

Avg score: 7.4/10

Sem Y+1:

Class size: 359

Avg score: 7.1/10



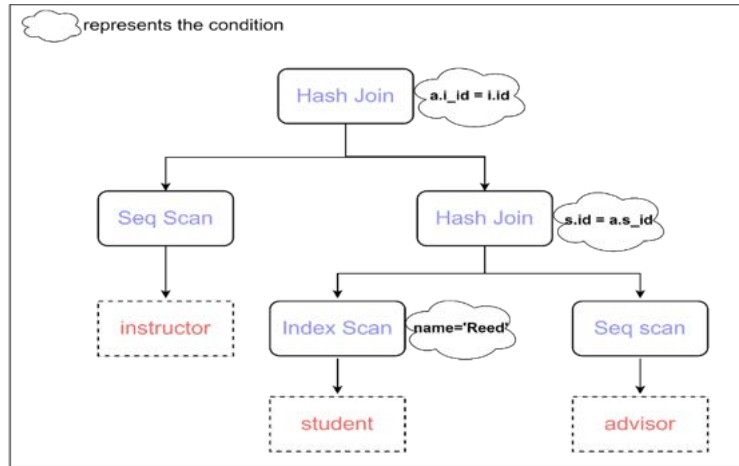
Is Understanding QEPs a Challenge to Learners?

Common Mistakes

- Incorrect ordering of steps
- Use relational algebra
- Writing SQL query
- Lumping several steps into single step
- Exclude filtering conditions in scan
- Missing intermediate relations
- Unclear specification of physical operators
- Attempt to describe implementation of various operators



Understanding Alternative Plan Choices Made By DBMS



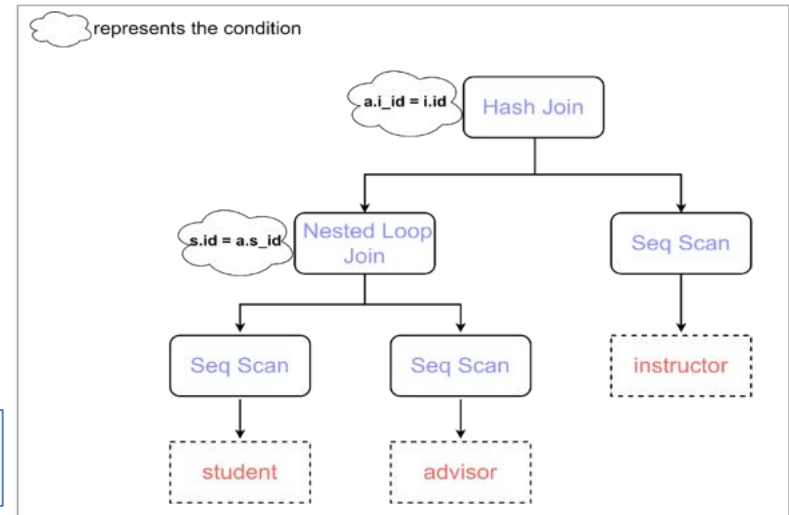
Plan A

Sem Z:

Class size: 188

Avg score: 5.1/15

Only 7.4% score 8 and above



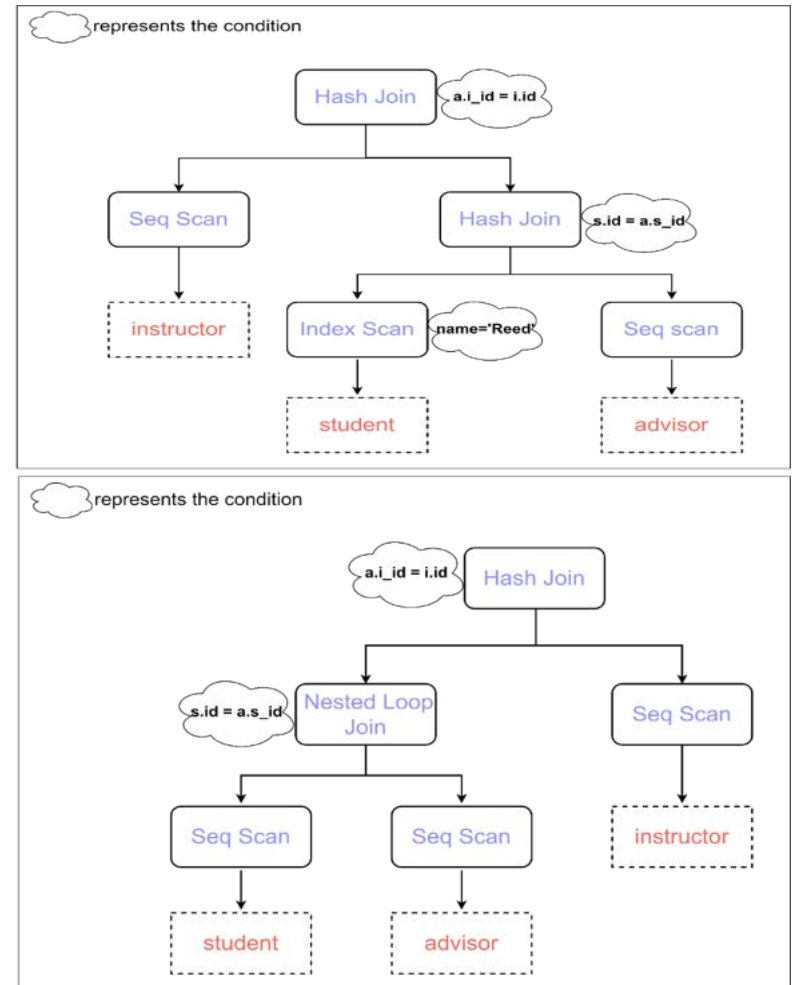
Plan B



Understanding Alternative Plan Choices Made by DBMS

Common Mistakes

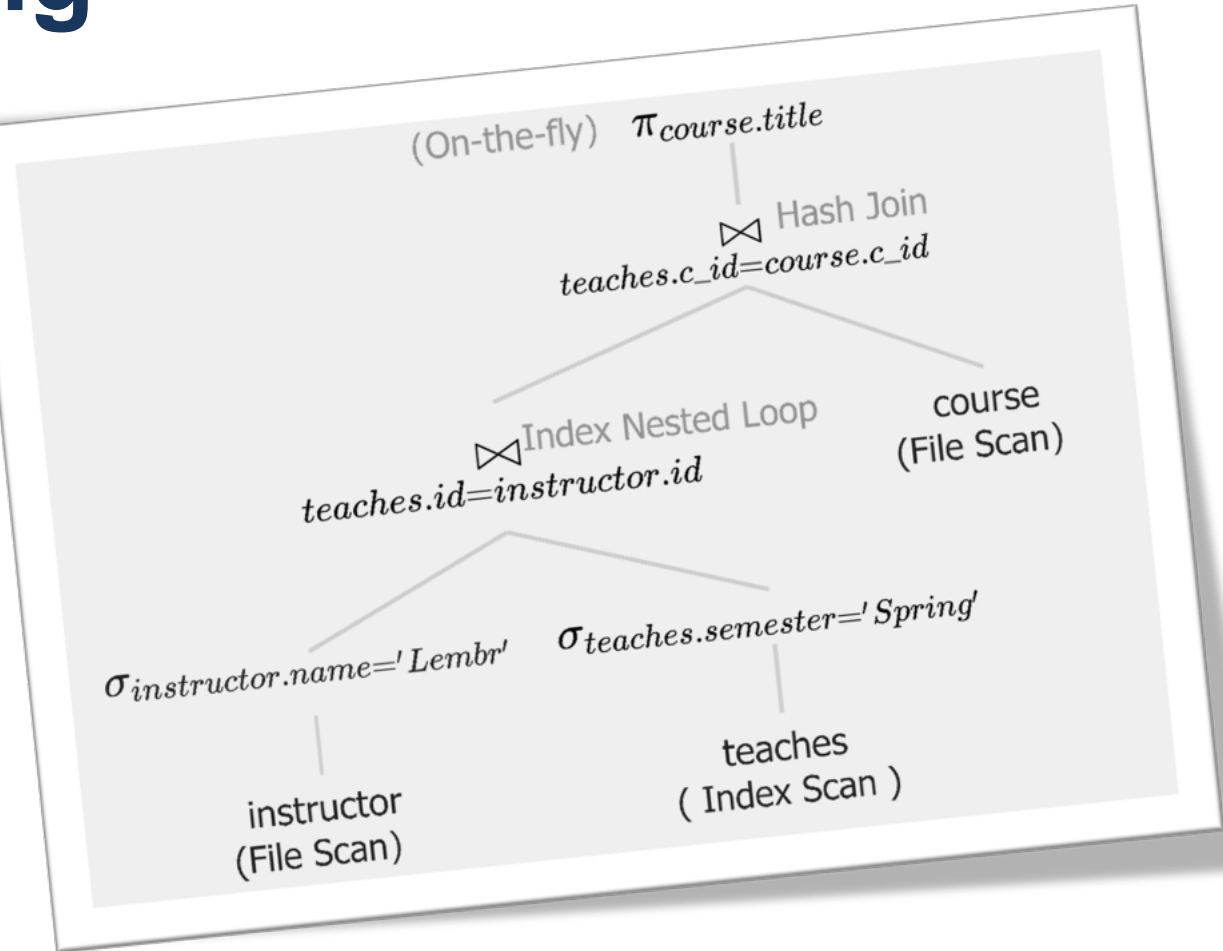
- Missing role of index scan
- Missing the possible impact of join ordering
- Explanation missing
- Incorrect justification w.r.t. the type of join operator



The Cost Estimation Challenge in Learning

Sem Y+1:

- 55% students scored less than 6/10
- 2/359 students got the cost estimation correct



The Cost Estimation Challenge In Learning

Common Mistakes

- Incorrect cardinality estimation of intermediate results
- Incorrect I/O cost of certain operations
- Inclusion of main memory cost

DBMS

Incorrect cardinality estimation

Deep learning-based techniques

Learners

Incorrect cardinality estimation

How to facilitate “deep” learning?



Limitations of Learning Modes

Textbook, lectures

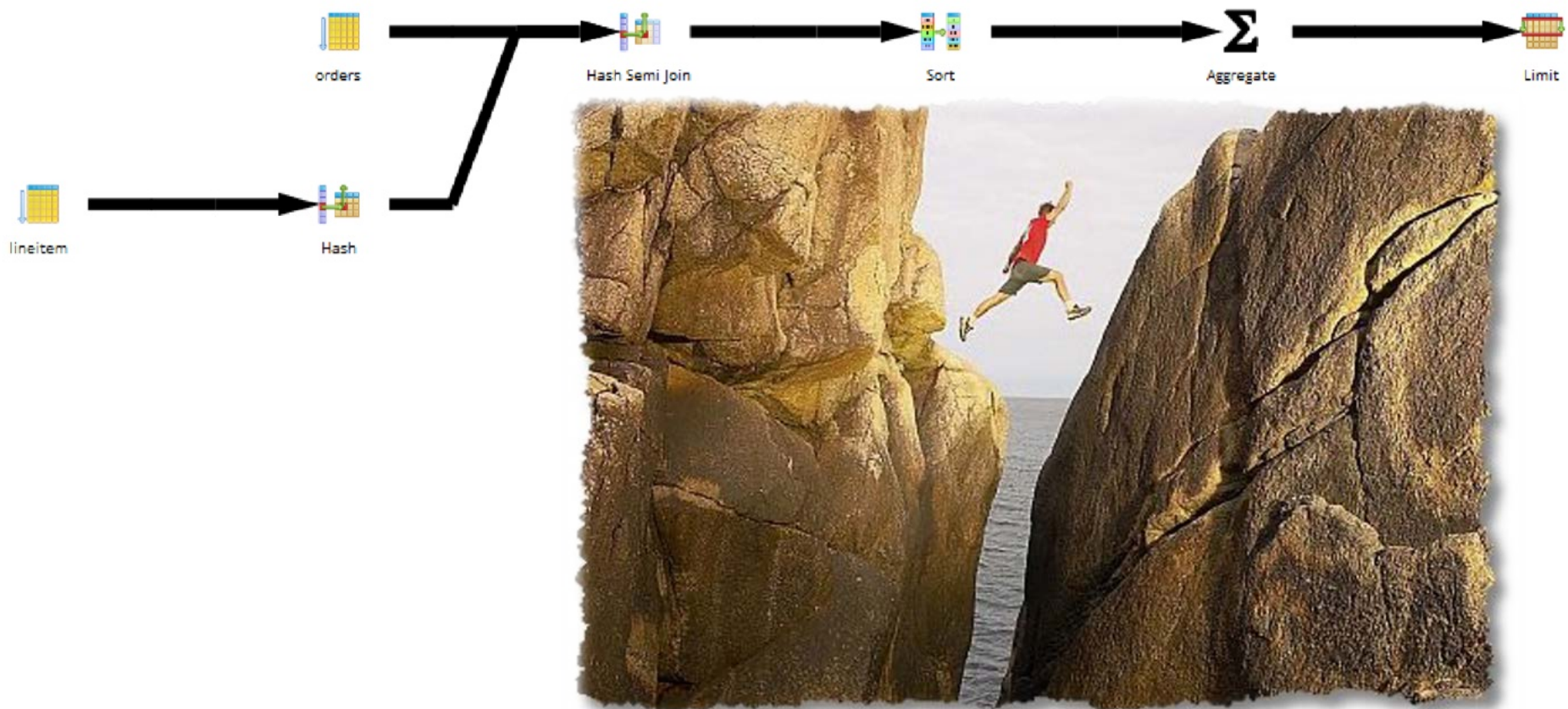
- Limited, hard-coded example problems
- Typically deals with simple SQL to illustrate concepts
- Not interactive
- You cannot learn about **any** SQL queries!

Off-the-shelf RDBMS

- For enterprise users
- Not designed for pedagogical support



All We Can Get from an RDBMS (Easily)....



Observations About Learners in Traditional Learning Environment

Learners in traditional settings

- Largely **extrinsically** motivated (e.g., getting good grades).
- **Learn-by-example**.
- Avoid exhaustive online search for resources and examples.
- Prefer slide decks and videos over textbooks.
- Limited by time due to concurrent courses.
- **Massing** vs spacing.



Question 2

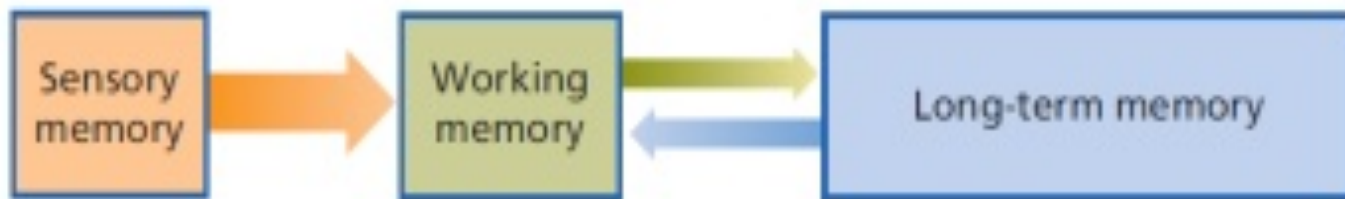
**Why learners face learning challenges
in traditional settings?**



Memory

Memory

- Takes meaningless sensory information (e.g., sound of professor's voice) as input
- Changes it into meaningful patterns (words, sentences, concepts) you can store and use later.



Memory is generally thought to be divided into three stages of processing

Basic Tasks of Memory

Encoding

- Memories for concepts usually require deliberate encoding effort to establish a usable memory
- Elaboration

Storage

- Retention of encoded material over time.

Retrieval

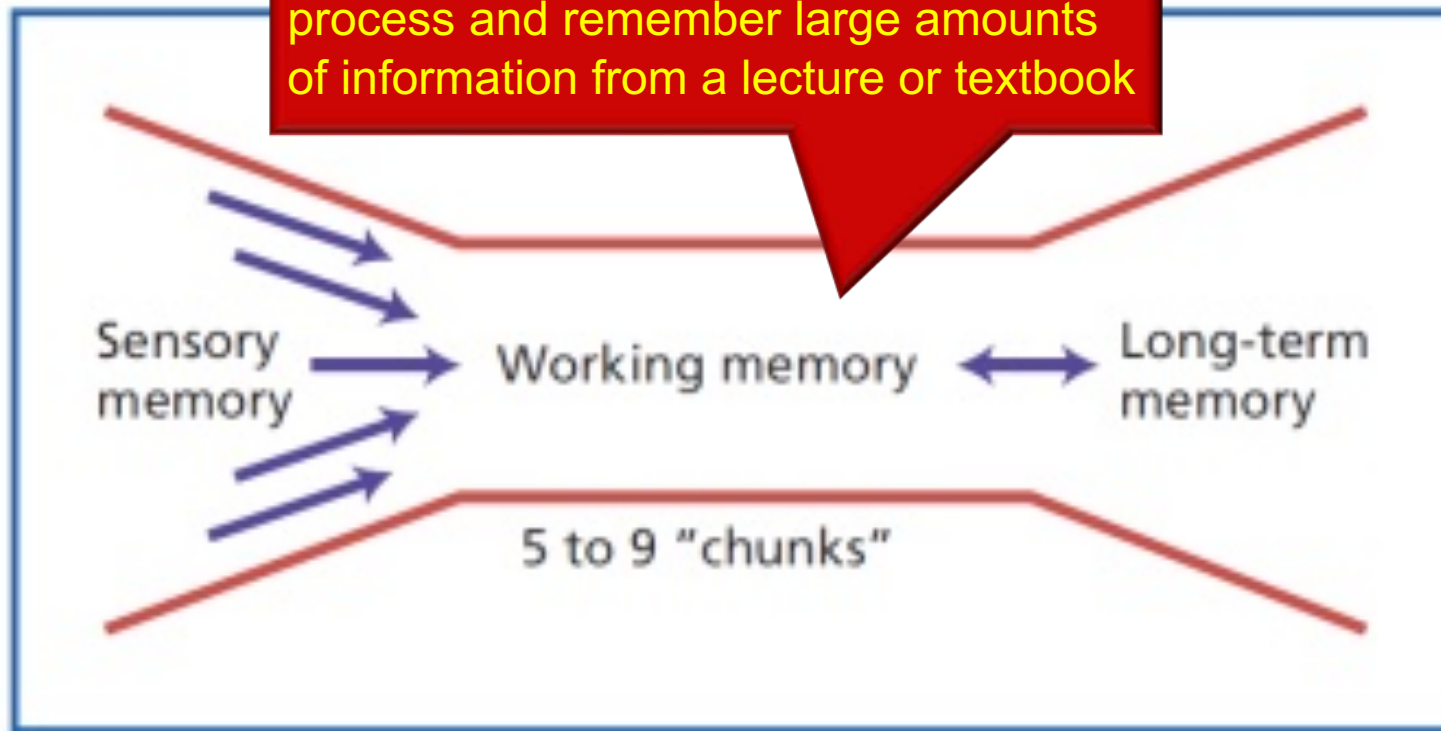
- Retrieve encoded memory accurately by exploiting good cue to access the information

Successful retrieval depends on how they were encoded and how they are cued.



How Do We Form Memory?

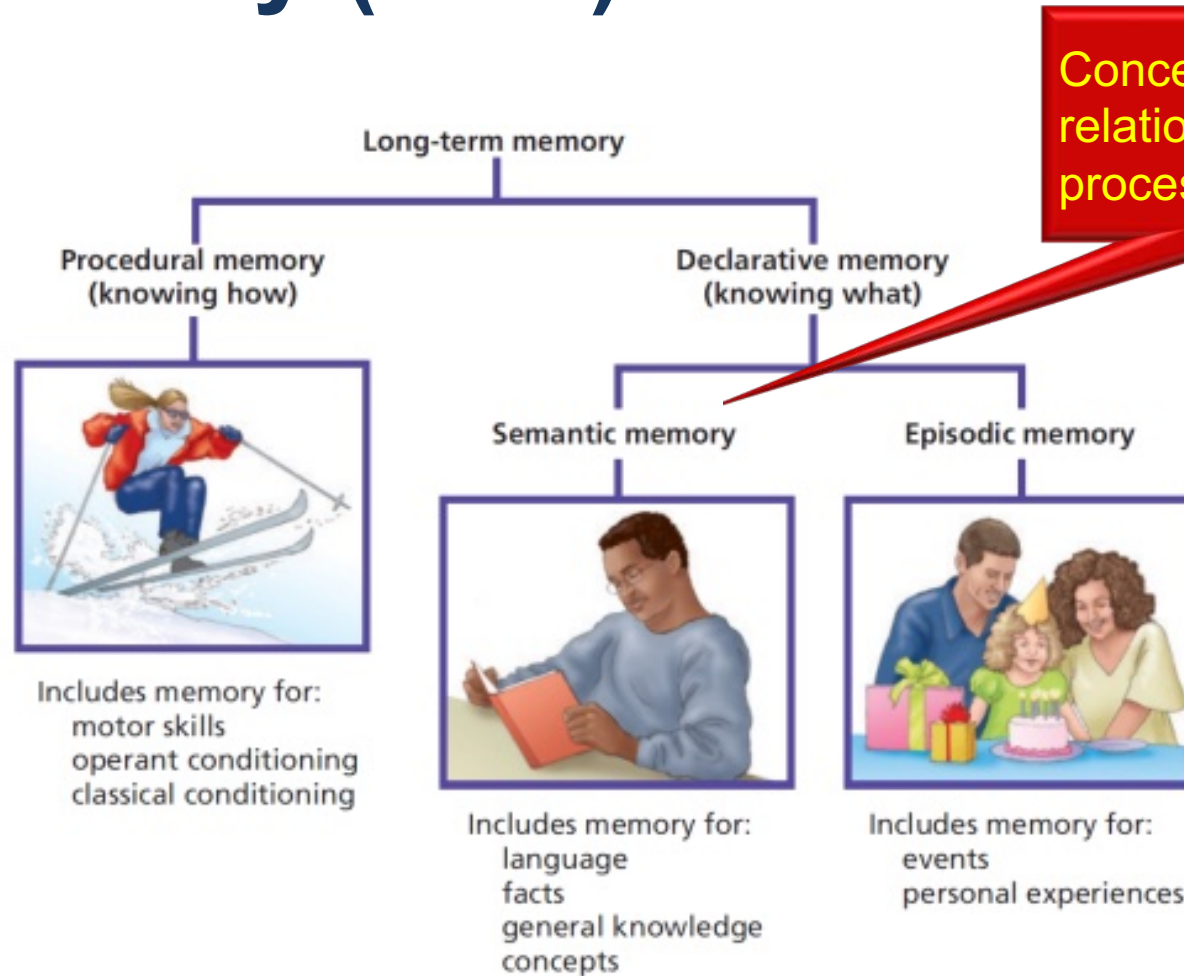
Major obstacle for students trying to process and remember large amounts of information from a lecture or textbook



P. Zimbardo, R. Johnson, V. McCann. Psychology Core Concepts. Pearson Education, Inc., 8th Edition, 2016.



Components of Long-Term Memory (LTM)

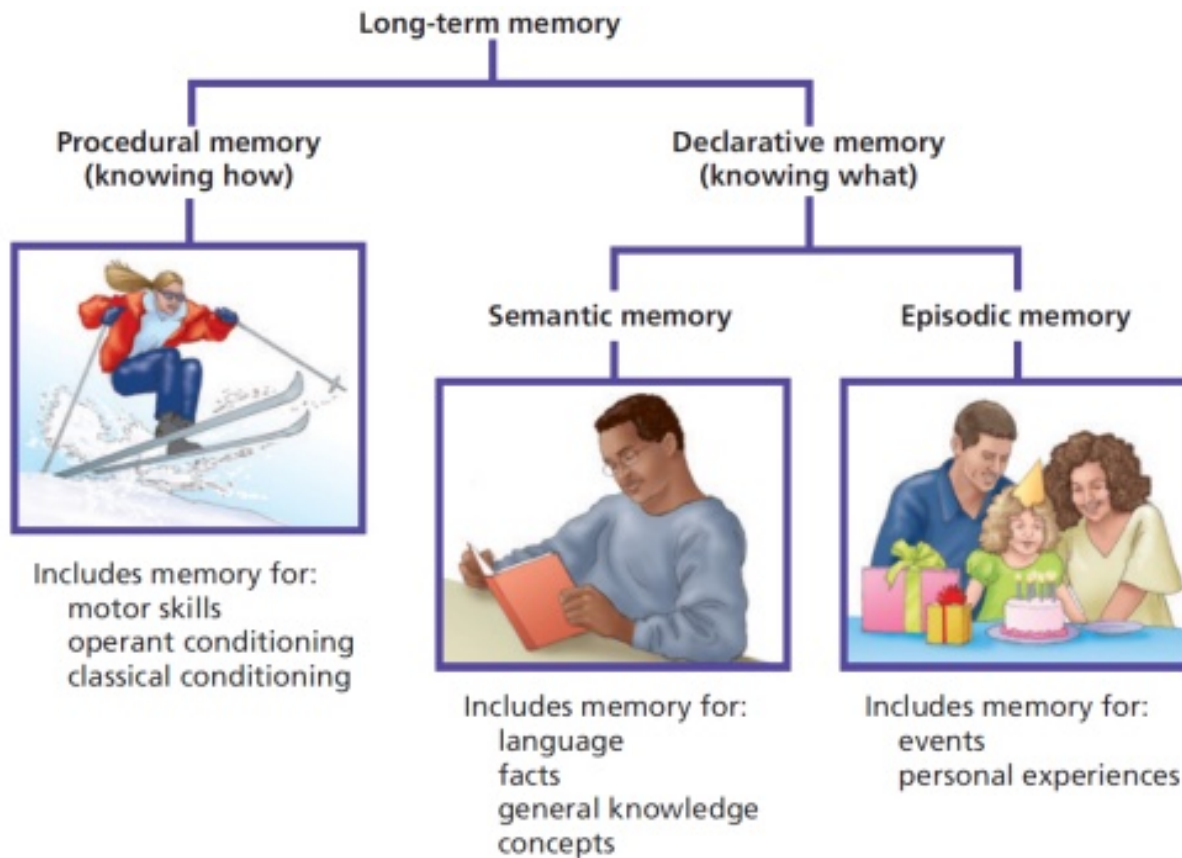


Concepts related to relational query processing

P. Zimbardo, R. Johnson, V. McCann. Psychology Core Concepts. Pearson Education, Inc., 8th Edition, 2016.



Biological Basis of LTM



Hippocampus and amygdala are crucial to laying down new declarative memories

- Memories gradually become more permanent with the help of the hippocampus.
- Memory consolidation.



“Seven Sins” of Memory

Sin	Description	Example
Transience	Decreasing accessibility of memory over time	Simple forgetting of long-past events
Absent-mindedness	Lapses of attention that result in forgetting	Forgetting location of car keys
Blocking	Information is present but temporarily accessible	Tip-of-the-tongue
Misattribution	Memories are attributed to an incorrect source	Confusing a dream for a memory
Suggestibility	Implanted memories about things that never occurred	Leading questions produce false memories
Bias	Current knowledge and beliefs distort our memories of the past	Recalling past attitudes in line with current attitudes
Persistence	Unwanted recollections that we can never forget	Traumatic war memories

P. Zimbardo, R. Johnson, V. McCann. Psychology Core Concepts. Pearson Education, Inc., 8th Edition, 2016.



Level-of-processing Theory

Lessons from Psychology

- The more connection you can make in working memory between new information and knowledge you already have, the more likely you are to remember it later

Level-of-processing Theory

- Craik and Lockhart (1972)
- “Deeper” processing establish more connections with LTM making new information more meaningful and memorable

P. Zimbardo, R. Johnson, V. McCann. Psychology Core Concepts. Pearson Education, Inc., 8th Edition, 2016.



How Do We Facilitate Deeper Processing?

Elaborative Rehearsal

- Putting concepts into your own words
- Adding examples that illustrate the concept – a type of elaborative rehearsal

Multi-modal interactions

- Multiple modes of interactions with the course material help to build a greater web of associations into which a memory is embedded.
- Learning theories: Learners learn better when the same content can be approached in multiple ways - both visual and verbal, as well as through hands-on learning.



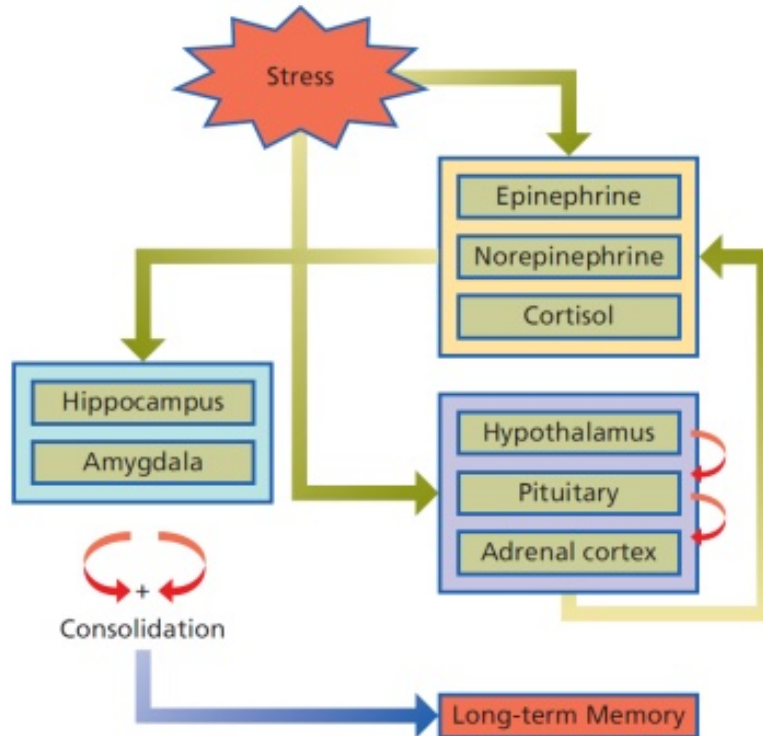
Memory Consolidation – Biological Basis

Biological Explanation

- New experiences consolidate much more rapidly (through hippocampus) if they are associated with existing **memory schemas**
- Why **elaborate rehearsal** and **depth of processing** help us to create more lasting memories



The Role of Stress



P. Zimbardo, R. Johnson, V. McCann. Psychology Core Concepts. Pearson Education, Inc., 8th Edition, 2016.

Release stress hormones in the brain, which act on the amygdala and hippocampus to strengthen the emotional memory of the event.

Remember how stressful DB course is instead of the concepts 😊 😊

Expectancy-Value Theory

If a task is too difficult or easy to complete then one may not engage with it.

Flow Theory

*A psychological state where a learner is **intrinsically** motivated to learn. Task is neither too easy or too difficult.*



Limitations of Traditional Modes

- Limited hard-coded examples

- Limited modes of approaching a content

- Processing a large amount of lecture and textbook content

- Stress due to massing, difficulty of accessing content



- Lesser web of association forming in long-term memory.



- Bottleneck of working memory.
- Impacts encoding and storage.



- Impacts memory consolidation in LTM.
- Expectancy Value Theory
- Flow Theory



Question 3: Towards Technology-Enabled Learning

**Can we build technologies to
supplement learning of relational
query processing?**



Broad Goals

Interactive

- Multi-modal, interactive mode
- Unlimited on-demand examples to facilitate elaborative rehearsal
- Create more web of connections

Easy

- Reduce difficulty in accessing information and content
- Facilitate **operant conditioning** through **negative reinforcement**
- Facilitate Expectancy value and flow theories
- Reduce stress

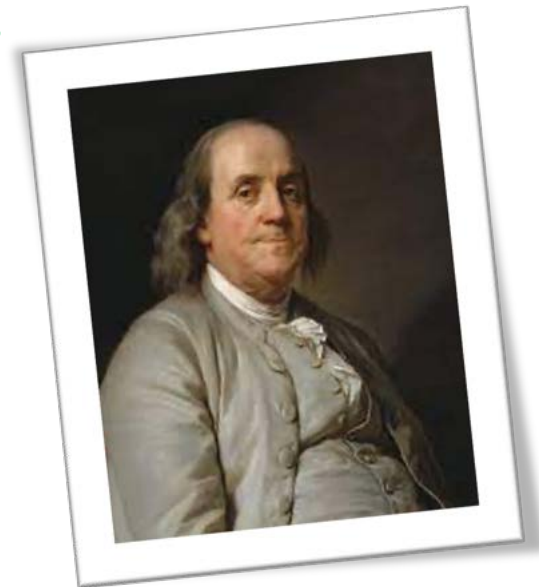
Understand

- Understand learning of learners through interaction data
- Feedback loop to improve pedagogy



Enhancing Learning Through Involvement

*“Tell Me and I Forget, Teach Me
and I May Remember,
Involve Me and I Learn.”*



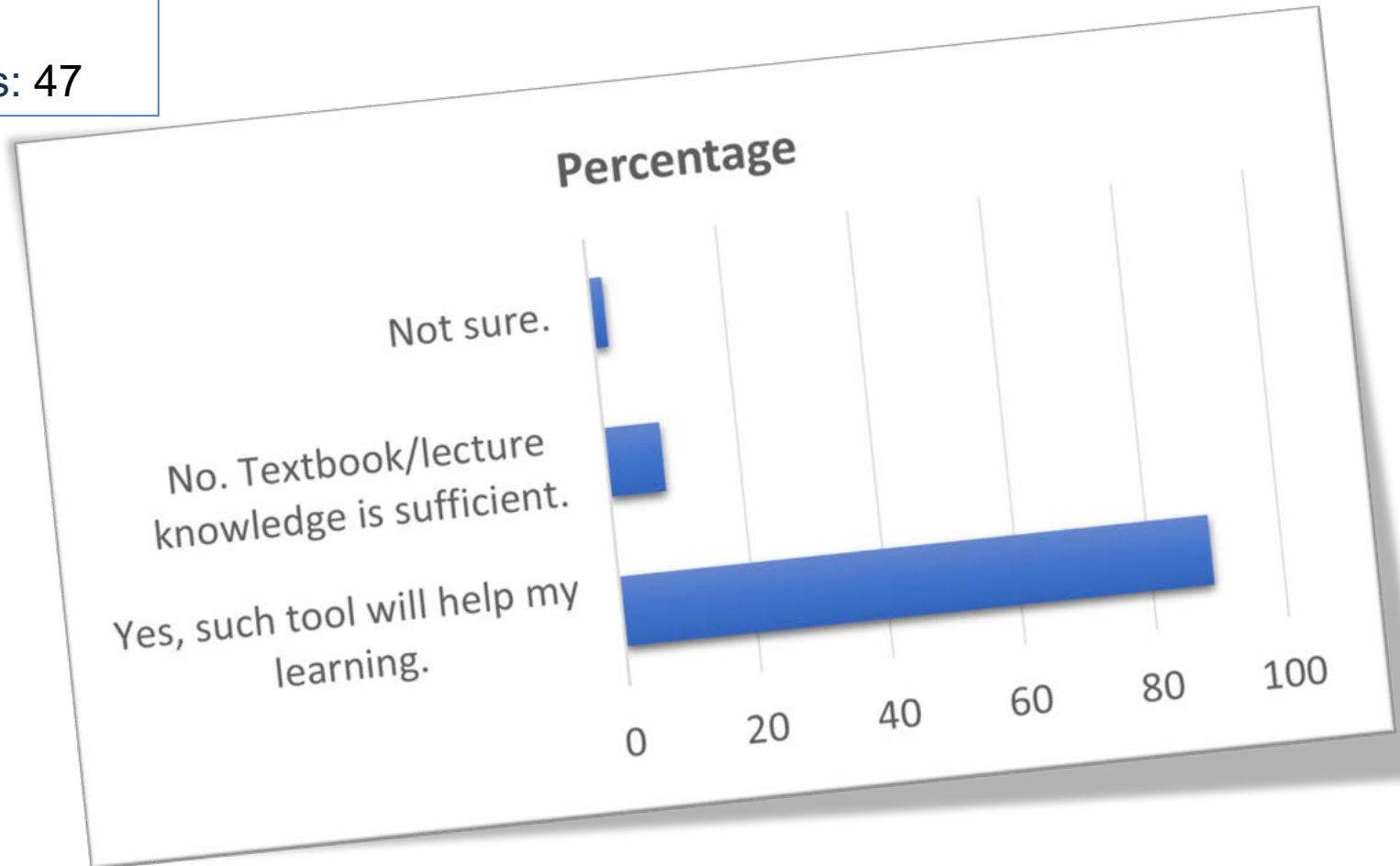
Benjamin Franklin



What Learners Think Of Tools to Augment Learning?

Sem X:

Respondents: 47



Understanding QEPs

Can we describe a QEP using visual and natural language to encourage involvement-based learning?



Challenges

Large QEP->NL training data is infeasible

Rule-based generation may create boredom

Should be generalizable w.r.t. application domain and RDBMS



NEURON & LANTERN

LANTERN

Switch DB POEM Menu ▾

Database schema
(current DB: TPCH)

▷ customer
▷ lineitem
▷ nation
▷ orders
▷ part
▷ partsupp
▷ region
▷ supplier

Example

```
select l_returnflag, l_line  
select l_orderkey, sum(l  
select o_orderpriority, c  
select n_name, sum(l_e  
select sum(l_extendedp
```

SQL query

```
select  
  l_orderkey,  
  sum(l_extendedprice * (1 -  
l_discount)) as revenue,  
  o_orderdate,  
  o_shippriority  
from  
  customer,  
  orders,  
  lineitem
```

Submit

Question

How many rows are left after a ce ▾

Step

Enter the step number

Submit

Natural language description of
QEP

The query is executed as follow.

Step 1, perform sequential scan on table lineitem and filtering on (l_extendedprice > '10') to get intermediate table T1 .

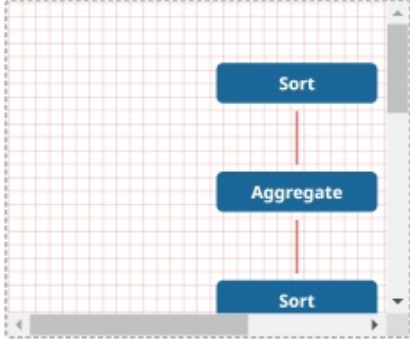
Step 2, perform sequential scan on table orders and filtering on (o_totalprice > '10') to get intermediate table T2 .

Step 3, perform sequential scan on

Compare

Answer

Visualize Plan



Detailed view

Feedback

Please provide your feedback here...

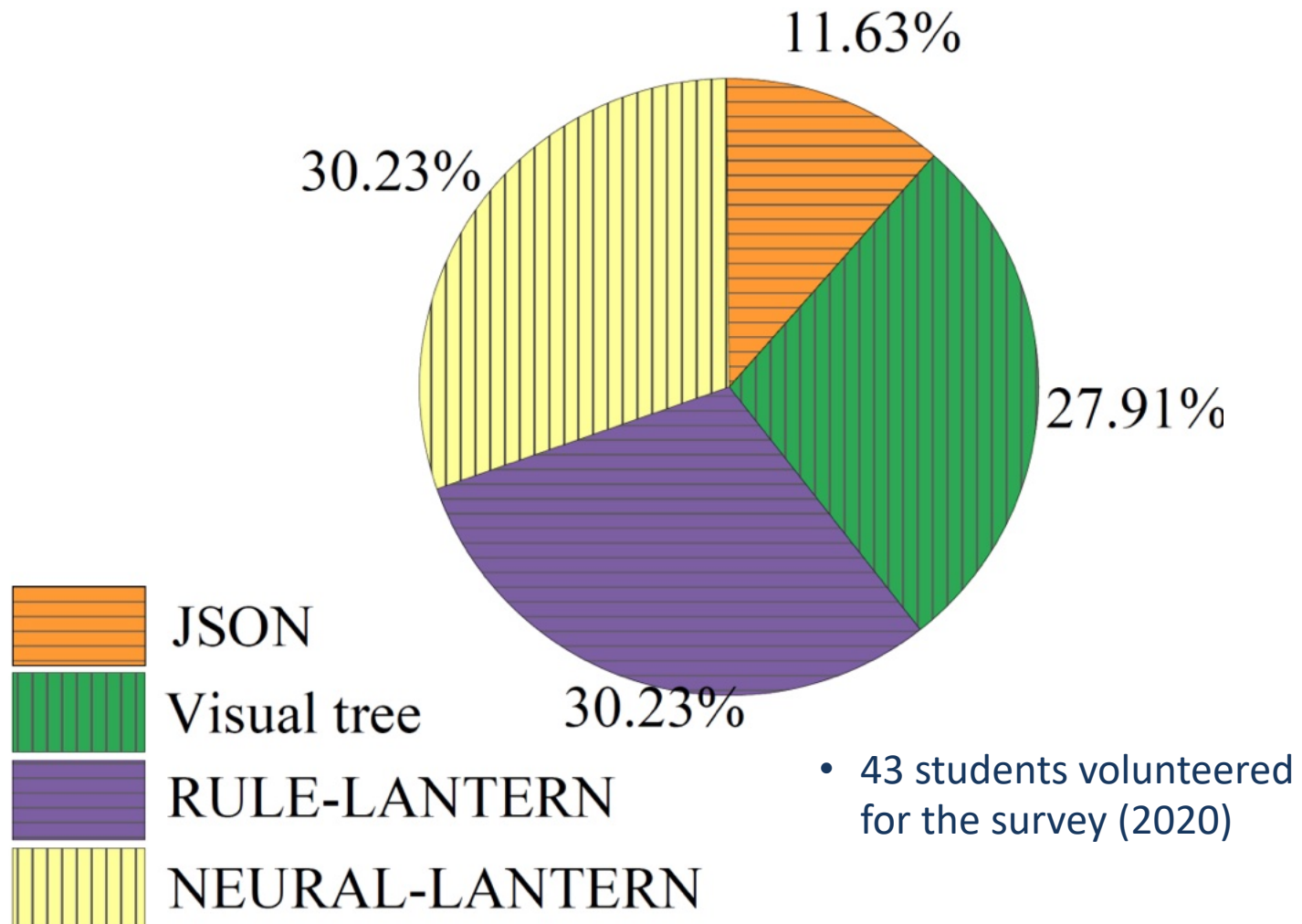
Submit

<https://howardlee.cn/lantern/>

- **NEURON: Query Execution Plan Meets Natural Language Processing For Augmenting DB Education.** Siyuan Liu, Sourav S Bhowmick, Wanlu Zhang, Shu Wang, Wanyi Huang, Shafiq Joty. In SIGMOD, 2019
- **Towards Enhancing Database Education: Natural Language Generation Meets Query Execution Plans.** Weiguo Wang, Sourav S Bhowmick, Hui Li, Siyuan Li, Shafiq Joty, Peng Chen. In SIGMOD, 2021.
- **LANTERN: Boredom-conscious Natural Language Description Generation of Query Execution Plans for Database Education.** Peng Chen, Hui Li, Sourav S Bhowmick, Shafiq R Joty, Weiguo Wang. In SIGMOD, 2022.



User Feedback: Which query plan format is most preferred?

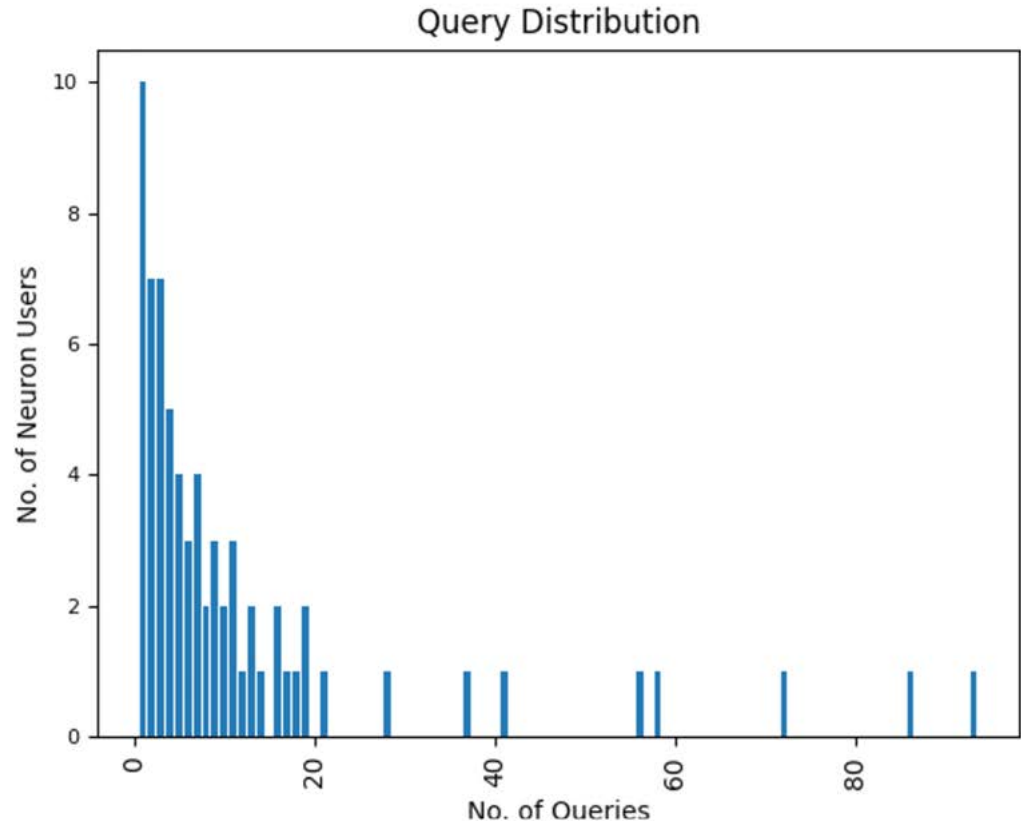


Impact of NEURON: When the Rubber Meets The Road

Sem X:

Class size: 166

- 41.5% used it
- No. of queries vs the number of distinct learners who posed that number of queries
- More than 85% of them posed more than one query



Sourav S Bhowmick, Hui Li. Towards Technology-Enabled Learning of Relational Query Processing. IEEE Data Engineering Bulletin, IEEE CS, September 2022



Test Performance

No. of students: 162

NEURON
users



Avg: 8.43

Max: 10

Min: 6.5

Median: 8

NEURON
non-users



Avg: 7.07

Max: 10

Min: 0

Median: 7.5



What-If Queries on QEPs

Students' Questions

- What is the impact on cost if operator A (e.g., hash join) is replaced by operator B (e.g., nested-loop join)?
- What will be the impact on cost if a specific join ordering is changed?
- What if the plan uses the index/sequential scan operator?



MOCHA

<https://howardlee.cn/mocha/>

MOCHA

Switch DB Menu

Database schema
(current DB: IMDB)

▷ actors

▷ directors_genres

▷ directors

▷ movies

▷ roles

▷ movies_directors

▷ movies_genres

Example

```
SELECT * FROM actors WHERE gender = 'M';
SELECT count(*) FROM directors; ...
SELECT role FROM roles, actors WHERE role = 'Director';
SELECT count(*) FROM directors_genres WHERE genre = 'Action';
SELECT count(*) FROM movies WHERE year = 2010;
SELECT count(*) FROM directors, movies WHERE director_id = movie_id;
```

SQL query

```
SELECT
  count(*)
FROM
  directors,
  movies_directors
WHERE
  directors.id = movies_directors.director_id
GROUP BY
  directors.id
ORDER BY
  COUNT(*) DESC;
```

Config Alternative Plans

Query

Max config parameters

6(default)

Select an approach

SINGLE plan(default)

Select parameter(s)

☐ Bitmap Scan

☒ Index Scan

☐ Index-Only Scan

☐ Sequential Scan

☐ TID Scan

☒ Hash Join

☐ Merge Join

☐ Nested-Loop Join

☐ Hashed Aggregation

☐ Materialization

☐ Explicit Sort

Visualize Plan

Choose one plan to view

AP1

Explanation

The QEP (SP) is selected because it has the least cost amongst other plans.

Merge joins are preferred if the join inputs are large and are sorted on their join column.

Total Cost

Plan	Total Cost
SP	29,036.38
AP1	20,000,074,289.69

SP	AP1
29,036.38	20,000,074,289.69

Detailed view

MOCHA: A Tool for Visualizing Impact of Operator Choices in Query Execution Plans for Database Education. Jess Tan, Desmond Yeo, Rachael Neoh, Huey Eng Chua, Sourav S Bhowmick. In VLDB, 2022.



Exploring Alternative Query Plan Space

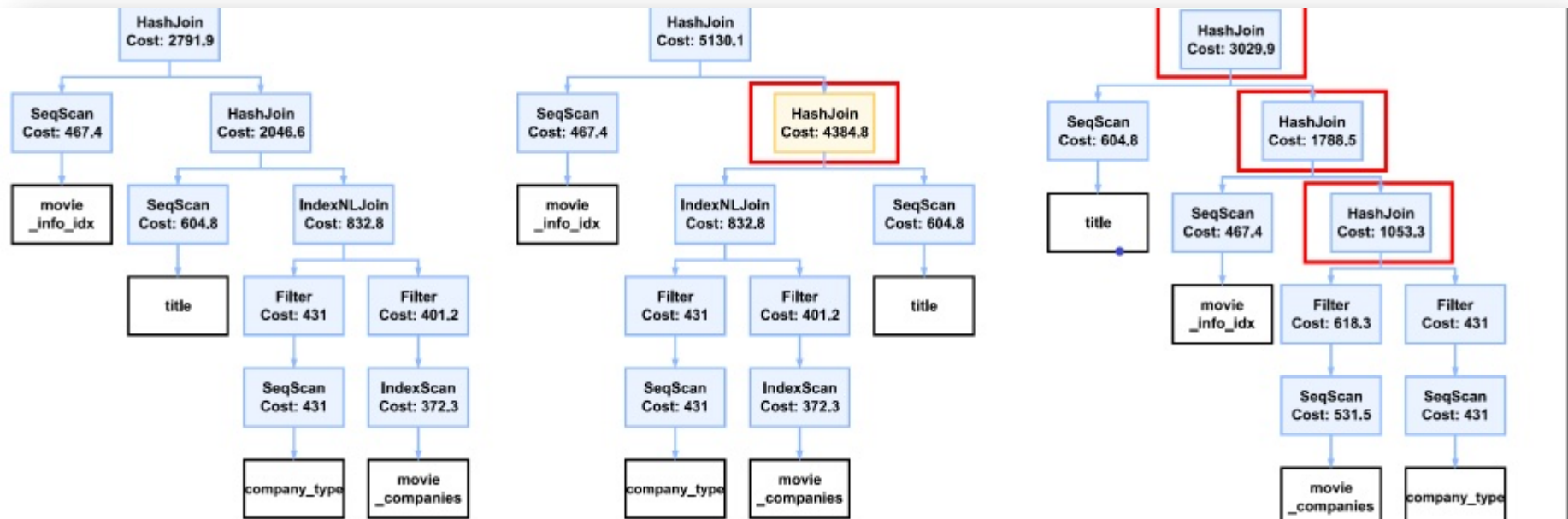
Can we explore alternative query plans
in a learner-friendly manner ?



Alternative Query Plans (AQPs)

Alternative Query Plans

- Given an SQL query, there are many different query plans, other than the QEP, for executing it.
- Alternative** query plans (AQP)



Challenges

Which plans are informative to learners?

How do we compute plan informativeness?

How do we design efficient algorithms?



ARENA

← → ↻ Not secure | 47.93.81.163/#/ ⌵ ☆ 📱 ⚙️ 🖨️ 🔄 Paused

ARENA

Database Information

Method

Parameter

Model B-AQP ☐ I-AQP

Tree Edit dist ☐ **# of plans**

Number

S_weight

C_weight

Cost_weight

lambda

Filter Limit

Setting

-- Query all courses selected by Reed and sort
-- 查询 Reed 选择的所有课程并排序

select title
from course, takes, student
where course.course_id = takes.course_id and
takes.id = student.id and student.name =
'Reed' order by title;

1.sql Execute

Id	S_dist	C_dist	Cost
0	0.00	0.00	991.3
1	0.37	0.32	9405.
2	0.00	0.27	7575.

```
graph TD; Sort[Sort Cost:7575.8] --> I1[InnerHashJoin Cost:7506...]; I1 --> T1[TableScan Cost:431.9]; I1 --> I2[InnerHashJoin Cost:7054...]; T1 --> course[course]; I2 --> T2[TableScan Cost:1262.6]; I2 --> F1[Filter Cost:431.9]; T2 --> takes[takes]; T2 --> T3[TableScan Cost:431.5]; T3 --> student[student];
```

全屏

放大 60%

缩小

刷新

下载

QEP

AQPs

ARENA: Alternative Relational Query Plan Exploration for Database Education. Hu Wang, Hui Li, Sourav S Bhowmick, Baochao Xu. In SIGMOD, June 2023

NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

Test Performance

No. of students: 50

ARENA users
(Gp 2)



Avg: 8.24

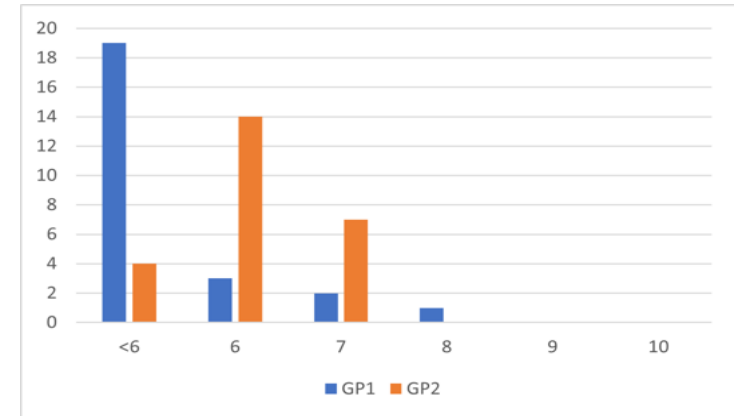
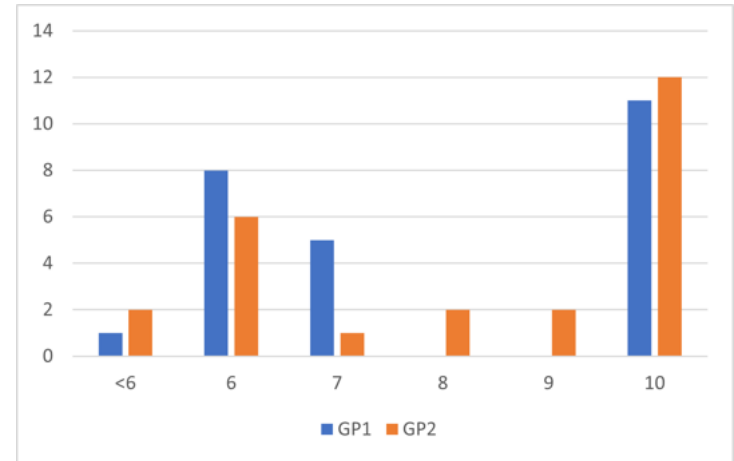
Avg: 6.04

ARENA non-
users (Gp 1)

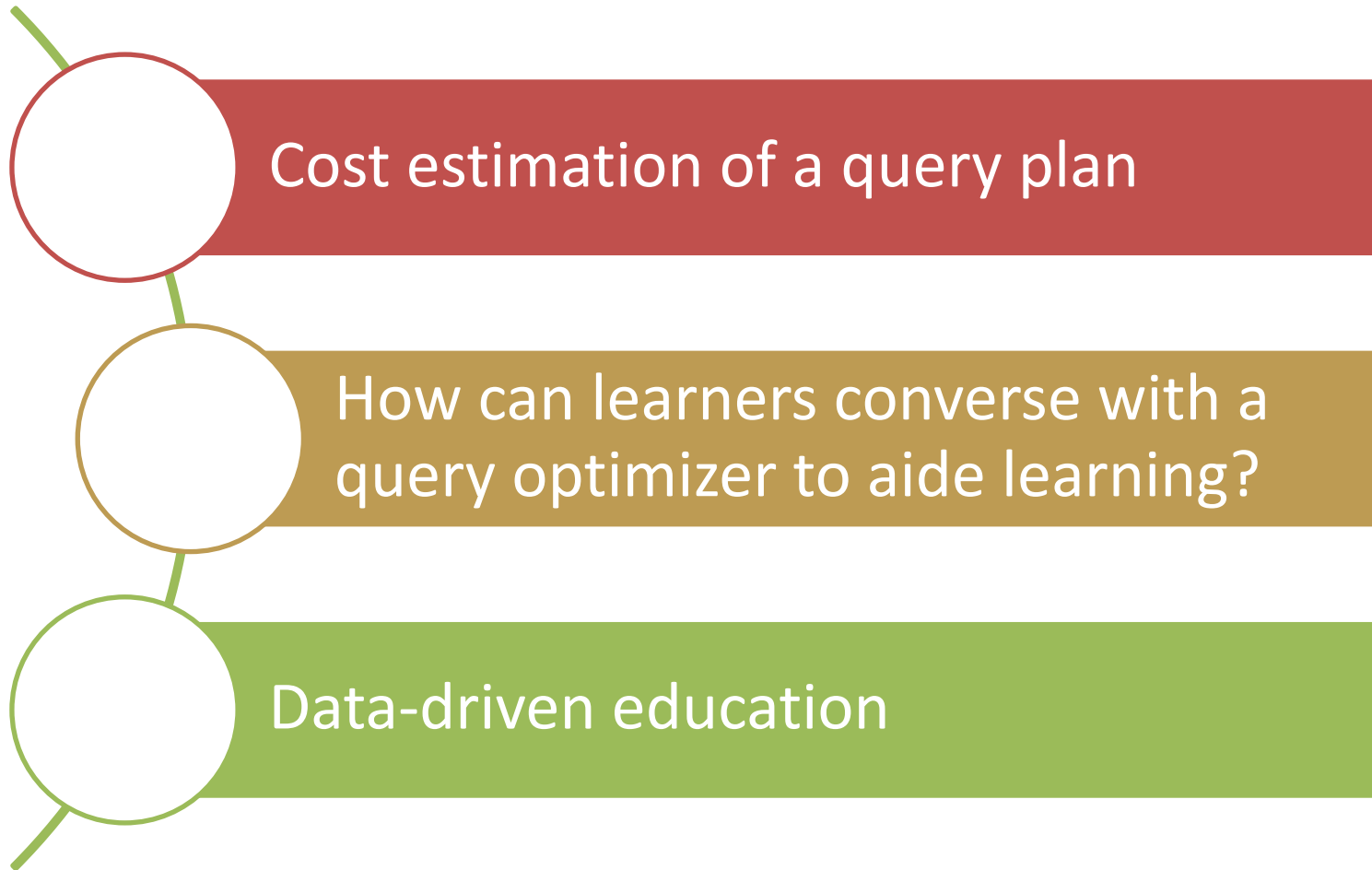


Avg: 7.76

Avg: 4.12



What's Next?



Learning Cost Estimation - Challenges

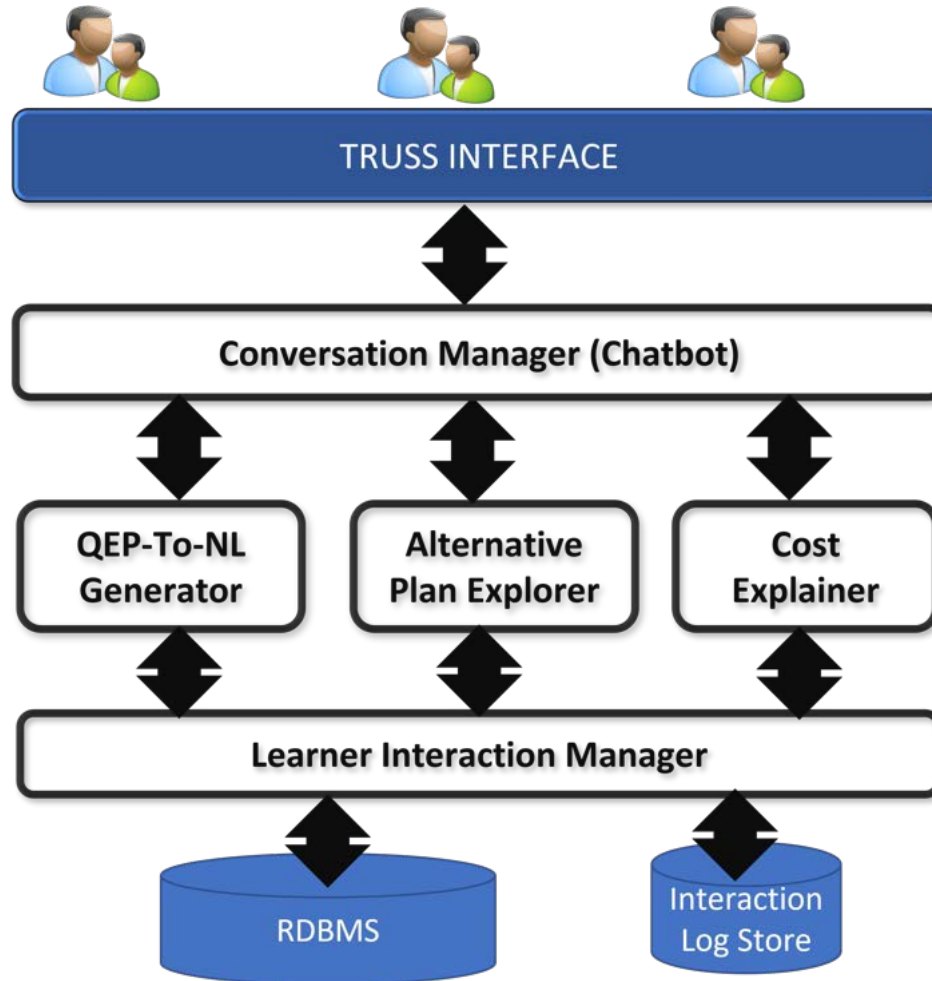
Extract cost computation formula for subtrees

Connecting with textbook knowledge

Explainable and generalized approach



TRUSS System



Sourav S Bhowmick, Hui Li. Towards Technology-Enabled Learning of Relational Query Processing. IEEE Data Engineering Bulletin, IEEE CS, September 2022



Conversation with TrussBot

TrussBot

Can you please explain how my query is executed?

<Describes the QEP in NL>

There is a zigzag join operator it seems in the plan. I do not think I have learnt it in the course. What is it?

<Defines zigzag join>

I see, How is the cost of this join estimated?

<Describes the formulas and steps for estimating the cost>

These formulas are from any textbook? Can you please explain the differences, if any?

<Explains in NL>

Why this plan is chosen? Can you connect the reasons with what we have learnt in textbooks?

<Explains in NL>

TrussBot

What will the plan be if the zigzag join is replaced by a hash join? Has such plan been considered by the optimizer?

Yes, it was considered. <Visualizes the plan with hash join and corresponding estimated cost>

Interesting! Can you show some alternative query plans considered by the engine that are interesting?

Let me show you one at a time. <Shows one interesting AQP at a time>.

Ok! I have seen a sufficient number of them. I think I have a better understanding of the different plan choices. Let me execute the query.

Sure! Go ahead.

Hmm! It is taking quite a long time to finish execution. Is an efficient plan chosen?

No, unfortunately a bad plan is chosen this time. <Explains the actual and estimated cost difference and reasons for slow execution>.



Towards Data-driven Education

Interaction Log

- Access time, duration, queries, interactions
- Analyze and correlate with academic outcomes

Insights

- What challenges they are facing?
- Learning preferences
- Massing vs spacing

Effectiveness

- Elaborate rehearsal and level-of-processing theory
- Any correlation between engagement of a platform and performance?



Do We Care About Disabled Learners?

DBMS for Whom?

- Primarily for able-bodied end users
- No systematic research on designing data management products for **disabled** users

DEI Matters!

- **Diverse** learners in lifelong learning environment
- **How can we facilitate learning for disabled learners?**



<https://www.henkel.com/company/diversity-and-inclusion>



How Do They Learn?

ASD/ADHD

- Autism Spectrum Disorder (ASD) and Attention Deficit Hyperactivity Disorder (ADHD)
- Face cognitive and behavioral challenges
- Classical mode of learning may be inadequate
- **Visual thinkers and learners** are common among ASD

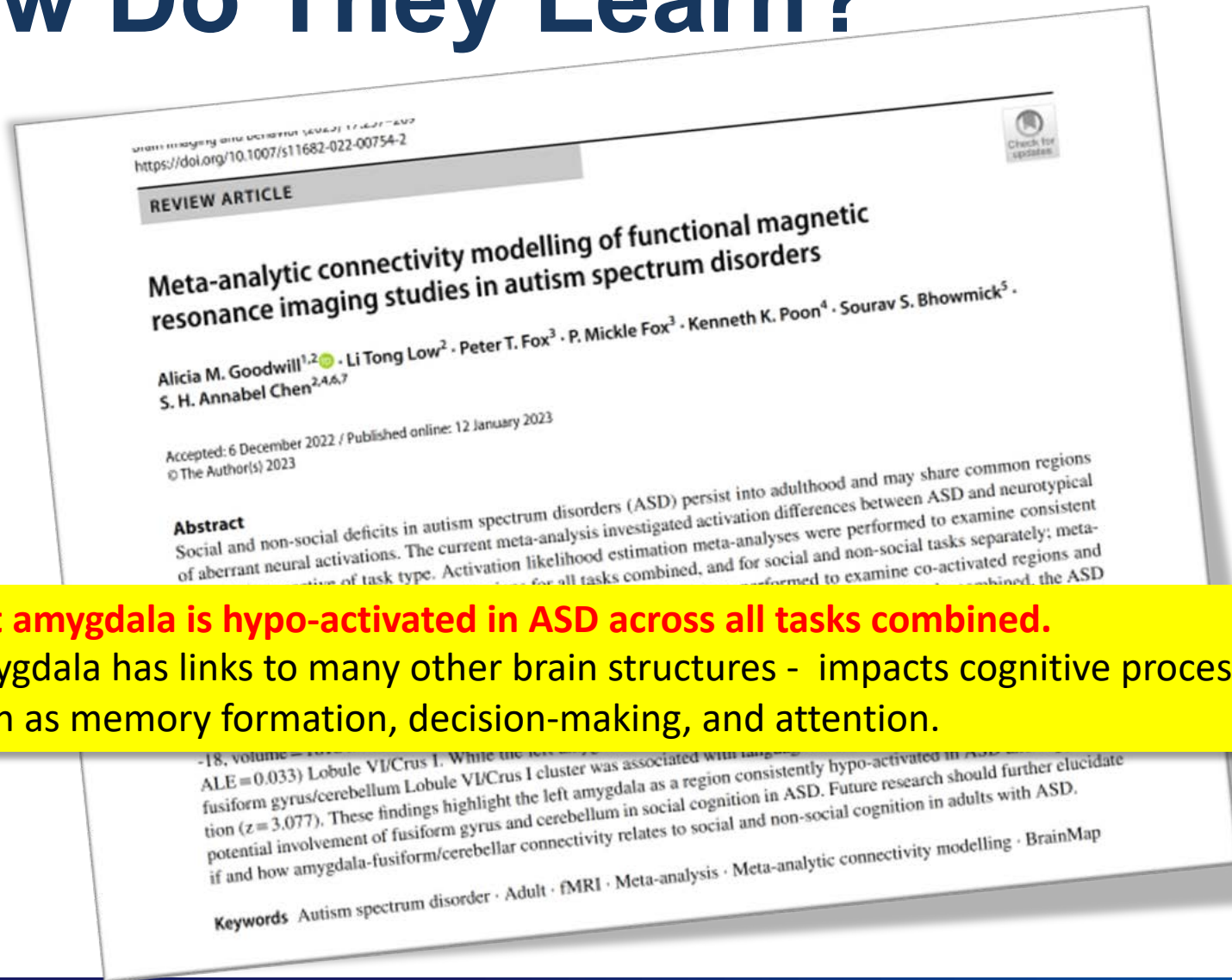
Understanding Learning

- **Differential** understanding of adult brain network and how it impacts learning.

Visual Tools That Can Help People With Autism Learn and Thrive. <https://www.verywellhealth.com/visual-thinking-and-autism-5119992>, 2021.



How Do They Learn?



- **Left amygdala is hypo-activated in ASD across all tasks combined.**
- Amygdala has links to many other brain structures - impacts cognitive processes such as memory formation, decision-making, and attention.

DB Researchers Needs to Break Out from the Enterprise Jar



How can technology supplement learning of database technology?


Beyond enterprise users


Implications to data science and data management education





Efforts on Query Visualization

Principles of Query Visualization

Wolfgang Gatterbauer 
Northeastern University
w.gatterbauer@northeastern.edu

Cody Dunne 
Northeastern University
c.dunne@northeastern.edu

H.V. Jagadish 
University of Michigan
jag@umich.edu

Mirek Riedewald 
Northeastern University
m.riedewald@northeastern.edu

Abstract

Query Visualization (QV) is the problem of transforming a given query into a graphical representation that helps humans understand its meaning. This task is notably different from designing a Visual Query Language (VQL) that helps a user compose a query. This article discusses the principles of relational query visualization and its potential for simplifying user interactions with relational data.

SQLVis: Visual Query Representations for Supporting SQL Learners

Daphne Mierdema, George Fletcher
Department of Mathematics and Computer Science
Eindhoven University of Technology
{d.c.mierdema, g.h.fletcher}@tue.nl

Abstract—SQL is a typical query language for performing data analytics. Although its usage is ubiquitous, learners experience that query formulation in SQL is error-prone and time-consuming. Prior research has shown that this is due to low expressive ease, extensive training requirements and high cognitive load, all of which present a significant burden for SQL learners. Visual representations can assist learners to significantly lower this burden. The current dominant paradigm aims to facilitate SQL querying by helping users to avoid the syntax of SQL. Such Visual Querying Systems (VQS), however, are not effective for SQL learners as they hide the syntax of the language during query formulation, rather than enabling learners to write correct queries in SQL. Furthermore, training with VQS is system specific, which leads to system dependency for learners. We argue that novices need support from Visual Query Representation (VQR) solutions which, instead, help them in learning how to write correct and portable SQL queries. In this paper we present SQLVis, a VQR to support novice SQL users in query writing. Our system represents the query as written in SQL by the user, which can improve the SQL writing proficiency of its users. Results of an in-depth empirical study demonstrate the significant value of SQLVis for learners.

Index Terms—Query languages, Visualization techniques and methodologies, Computer Science education



Fig. 1: A user query and its corresponding SQLVis visualization. Tables are encoded as nodes, and constraints that link tables together as edges. The returned attributes are highlighted in orange, column based constraints in green.

leads to increased information throughput [8]. However, they also increase system dependency, as a VQS obfuscates syntax instead of supporting the learning of SQL.

Another approach to using visualizations to support query formulation is to generate a representation from a written



Acknowledgements



Hui Li, Xidian Univ, China

- Siyuan Liu, NTU
- Weiguo Wang, Xidian
- Peng Chen, Xidian
- Zheng Li, Xidian
- Hu Wang, Xidian

LANTERN

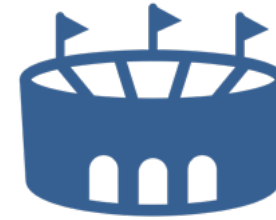



783



51

ARENA




688



35

* Since January 2022

