# LearnSphere to Integrate DataShop, MOOCdb, DataStage, DiscourseDB …
## Integrating Data Repositories Panel

Ken Koedinger

Professor of Human-Computer Interaction & Psychology

Carnegie Mellon University

Director of

**LearnLab**
Pittsburgh Science of Learning Center

Workshop 2: Advancing Data-Intensive Research in Education
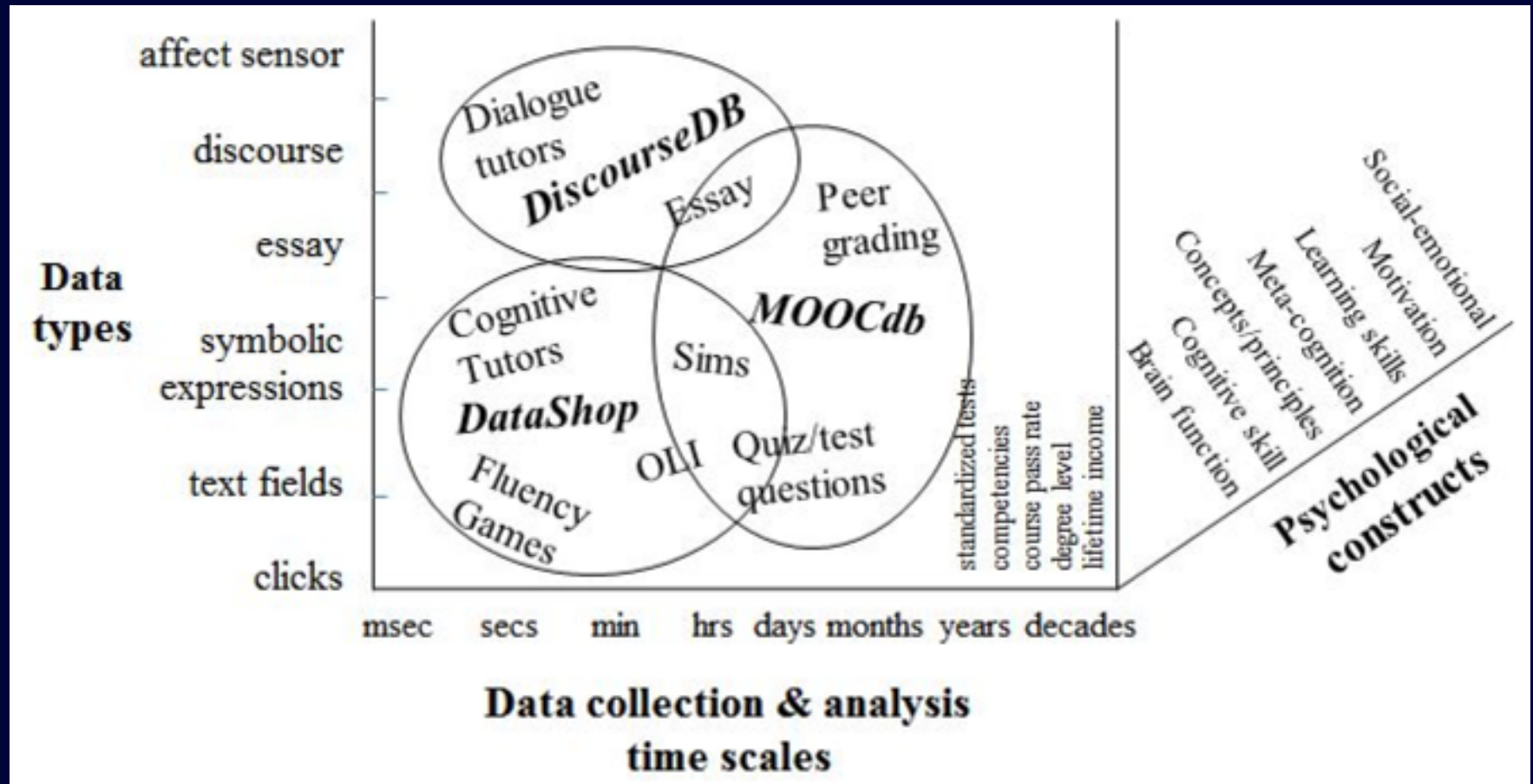June 1, 2015

# *Big Data* for education

## More important than "big"

- Collected as part of *natural* activities
- Affords experimentation, "A/B testing"

## Many dimensions of "big"

- *Tall* in number of participants (students)
- *Wide* in observations per participant (student)
- *Fine* in frequency of observation
- *Long* in spanning months or years
- *Deep* in theory-relevant variables

# LearnSphere: Integrate across data repositories toward answering questions



We need a education data infrastructures to integrate analytic methods
=> *produce discoveries not possible within current data silos*

# Cognitive Tutors
## Example source of educational data

My current cell phone company charges me $14.95 per month for service and $.13 per minute. PPS Cellular Phone Company has offered me $15.00 worth of free calls a month if I switch, but the charge is $.39 per minute.

| Quantity Name | Time | Current cost |
|---|---|---|
| Unit | minutes | $ |
| Expression | $t$ | .13t |
| Question 1 | | |
| Question 2 | | |
| Question 3 | | |
| Question 4 | | |

The cost from my current company increases by 0.13 each minute, but remember that it starts at 14.95 dollars.

*Fine*

*Wide*

**Authentic problems**

**Feedback *within* complex solutions**

Progress…

**Cognitive Tutor Algebra I**

### Scenario

My current cell phone company charges me $14.95 per month for service and $.13 per minute. PPS Cellular Phone Company has offered me $15.00 worth of free calls a month if I switch, but the charge is $.39 per minute.

**1.** How many minutes of calls can I get from PPS Cellular Phone Company for $50? What is the cost from my current company for that number of minutes?

**2.** How many minutes of calls can I get from my current company for fifty dollars? What is the cost from PPS Cellular Phone Company for that number of minutes?

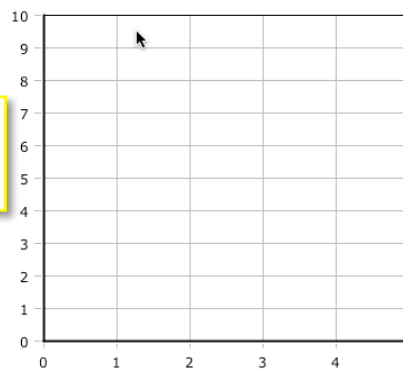**3.** What is the cost from both companies for sixty minutes

**Challenging questions**

...be the same?

**4.** After how many minutes of calls will the cost for both companies be the same?

### Worksheet

| Quantity Name | | | |
|---|---|---|---|
| Unit | | | |
| Expression | | | |
| Question 1 | | | |
| Question 2 | | | |
| Question 3 | | | |
| Question 4 | | | |

**Personalized instruction**

### Grapher

10.0

0.0

0.0

Legend: Enter Label

Equations: y = Enter Equation

**Hint**

If the cost from my current company and the cost from PPS Cellular Phone Company are equal, then their expressions are equal. Write an equation and solve it to find the number of minutes.

Close | << Previous Hint | Next Hint >>
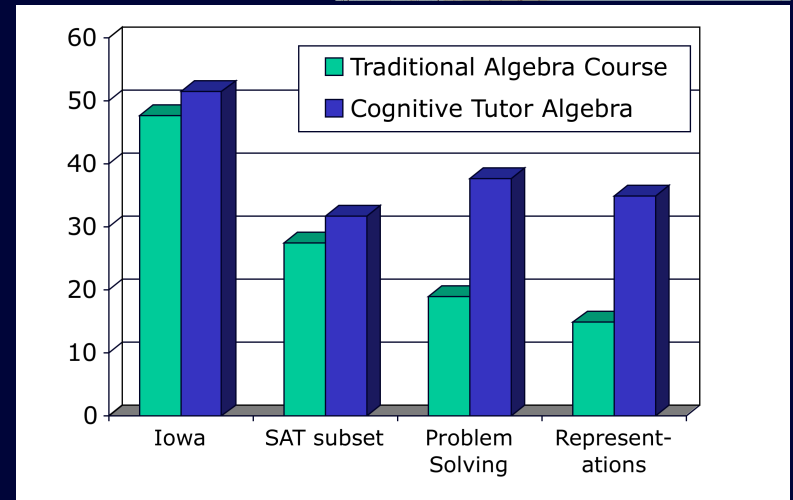
… **individualization**

Calculate input value.
Writing expression, any form.
Set axis bounds.
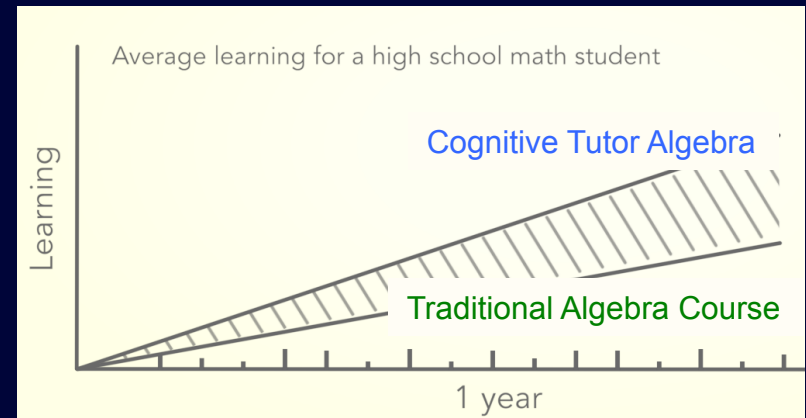Label point of intersection.
Enter given.

*Deep*

# Real World Impact of Cognitive Science

*Algebra Cognitive Tutor*

- Widespread intensive use
  ~600K students per year
  ~80 minutes per week

- *Many* field trials =>
  Student learning
  is 2x better

- Still:
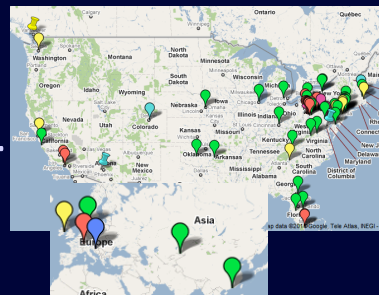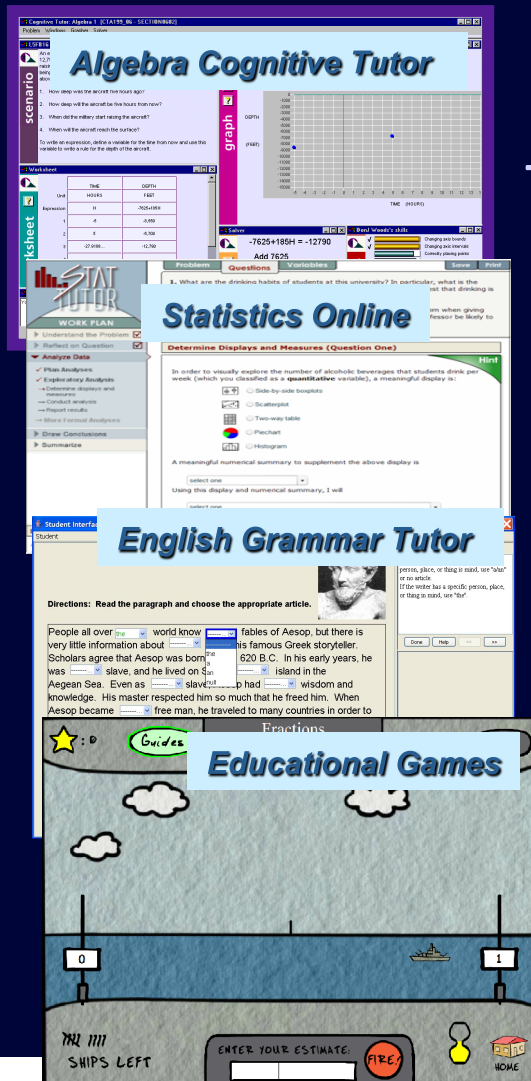  Could do better
  Too many decisions
  driven by intuition



Koedinger, Anderson, Hadley, & Mark (1997).
Intelligent tutoring goes to school in the big city.



Pane et al. (2013). Effectiveness of Cognitive Tutor
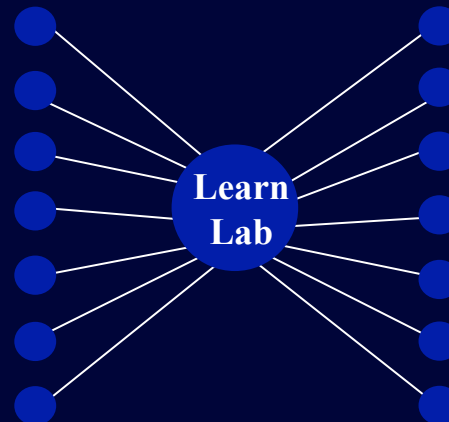Algebra I at Scale. Santa Monica, CA: RAND Corp.

# Social-technical infrastructure to discover conditions that cause *robust learning*

Ed tech    +    wide use    =    "Basic research *at scale*"



*Algebra Cognitive Tutor*

*Statistics Online*

*English Grammar Tutor*

*Educational Games*

+

=



**LearnLab**
Pittsburgh Science of Learning Center

Researchers                Schools

Learn Lab

Since 2004
> 680 ed tech data sets in DataShop

> 320 *in vivo* experiments

Koedinger et al. (2012).  The Knowledge-Learning-Instruction (KLI) framework:
Bridging the science-practice chasm to enhance robust student learning. *Cognitive Science*.

PSLC DataShop
a data analysis service for the learning science community

http://learnlab.org/datashop

**Help**

**Explore**
Public Datasets
Private Datasets
External Tools
What can I do?

**Learn More**
Documentation
About DataShop

# Welcome to DataShop, the world's largest repository of learning interaction data.

**Create an account** or **Log in** to start analyzing data.

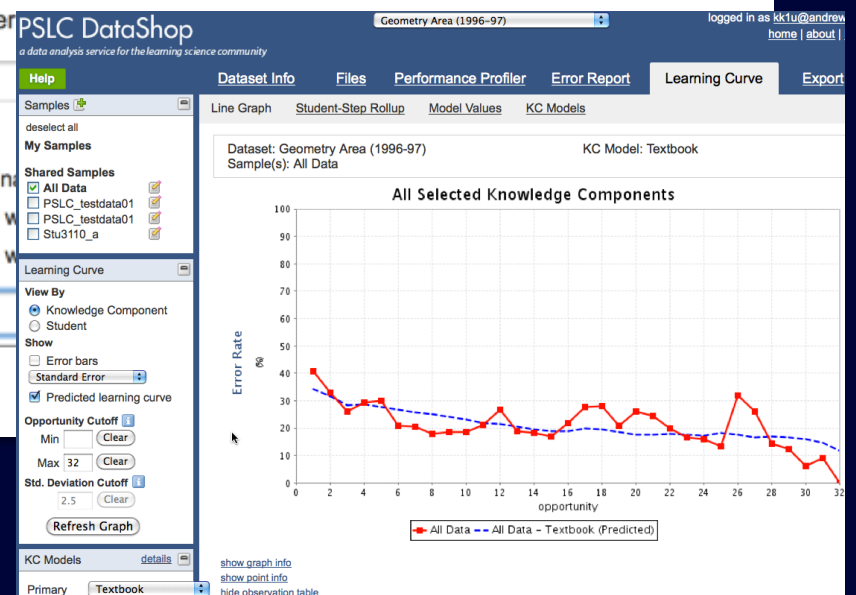## What can I do with DataShop?

## Upload a dataset

**Project** Add this dataset to ...
○ new project ○ existing project ○ choose later

**Project Name** Psychology MOOC data

**Data Collection Type**
○ Not specified
○ Not human subjects data (not origin...
● Study data collected under an IRB w
○ Study data collected under an IRB w

**Dataset Name** 2013 Psych

**Description** (optional)

PSLC DataShop
a data analysis service for the learning science community

Geometry Area (1996-97)

logged in as kk1u@andrew
home | about |

**Help**

Dataset Info | Files | Performance Profiler | Error Report | Learning Curve | Export

Line Graph | Student-Step Rollup | Model Values | KC Models

**Samples**
deselect all
**My Samples**
**Shared Samples**
☑ All Data
☐ PSLC_testdata01
☐ PSLC_testdata01
☐ Stu3110_a

Dataset: Geometry Area (1996-97)          KC Model: Textbook
Sample(s): All Data

**All Selected Knowledge Components**

**Learning Curve**
**View By**
● Knowledge Component
○ Student
**Show**
☐ Error bars
Standard Error
☑ Predicted learning curve
**Opportunity Cutoff**
Min          Clear
Max  32     Clear
**Std. Deviation Cutoff**
2.5          Clear
Refresh Graph

**KC Models**          details
Primary  Textbook

show graph info
show point info
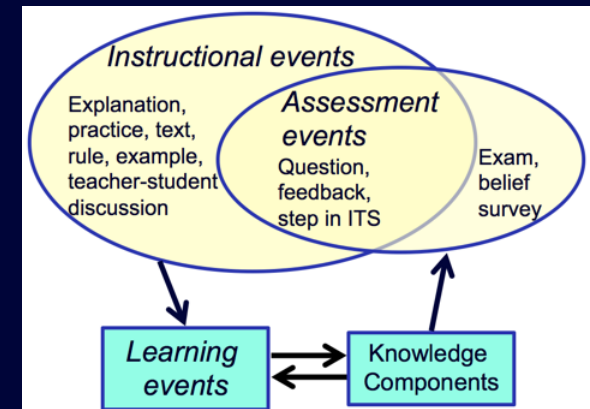hide observation table

— All Data --- All Data – Textbook (Predicted)

# 680 data sets
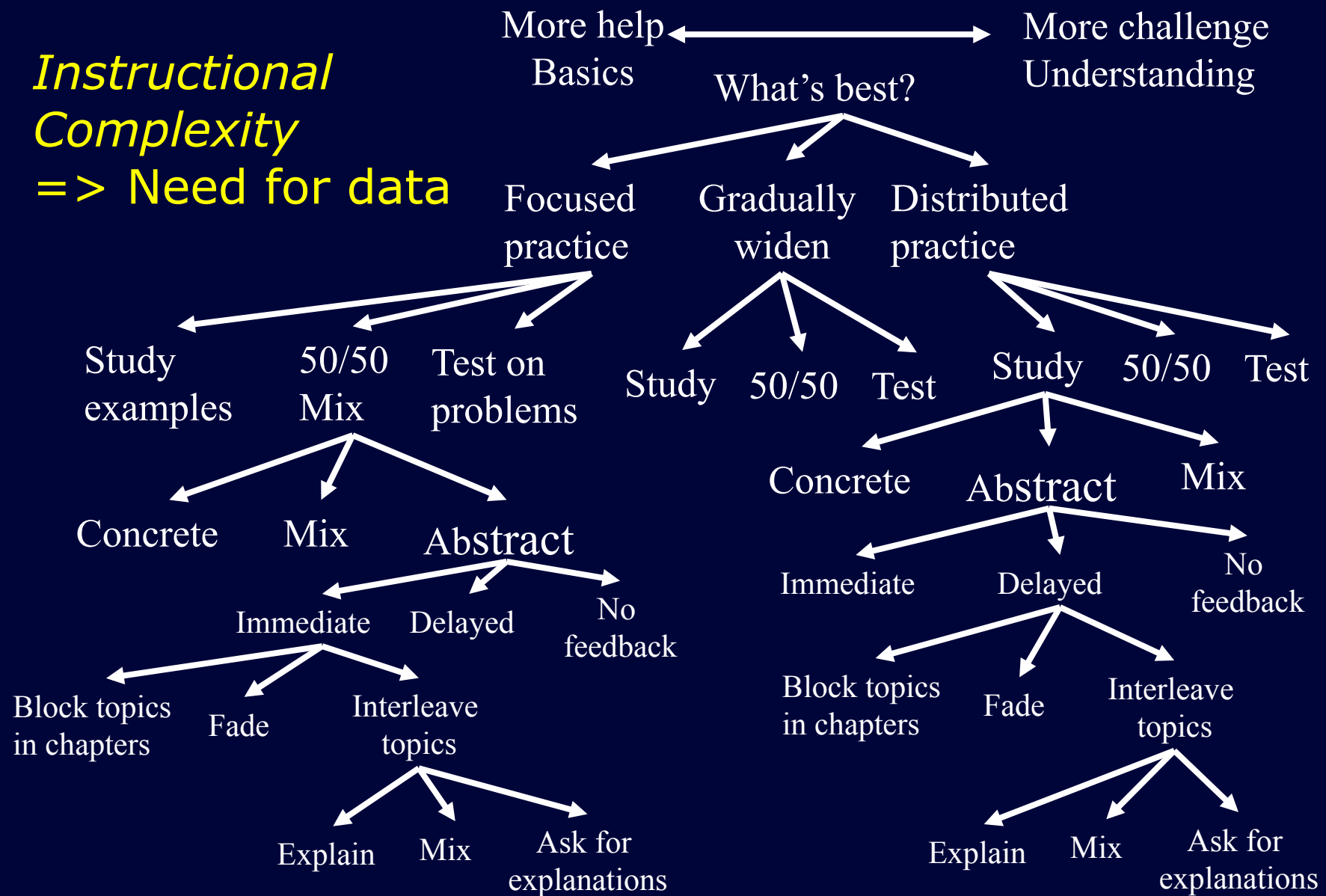## math, science, language …
## K12 & college

7

# Integrate across data repositories to answer questions

- Many complex *open questions* about the nature of:
  - Knowledge & cognition
  - Learning, metacognition
    - Motivation, & self-regulation
  - Instruction
- Need to work together to tackle these complex issues
  - Need to build on existing cognitive, social, education theory

Koedinger et al. (2012). The Knowledge-Learning-Instruction (KLI) framework. *Cognitive Science*.

8

*Instructional Complexity* => Need for data

More help ⟷ More challenge

Basics — Understanding

What's best?

Focused practice — Gradually widen — Distributed practice

**Focused practice:**
Study examples — 50/50 Mix — Test on problems

**Gradually widen:**
Study — 50/50 — Test

**Distributed practice:**
Study — 50/50 — Test

50/50 Mix: Concrete — Mix — Abstract

Abstract: Immediate — Delayed — No feedback

Immediate: Block topics in chapters — Fade — Interleave topics

Interleave topics: Explain — Mix — Ask for explanations

Study (distributed): Concrete — Abstract — Mix

Abstract: Immediate — Delayed — No feedback

Delayed: Block topics in chapters — Fade — Interleave topics

Interleave topics: Explain — Mix — Ask for explanations

Many other choices: animations vs. diagrams vs. not, audio vs. text vs. both, …

Koedinger, Booth, Klahr (2013). Instructional Complexity and the Science to Constrain It. *Science*.

$>3^{15*2} = 205$ trillion options!

9

# Automated support for cognitive task analysis: Discovering *hidden skills* using educational data

Cen, H., Koedinger, K., Junker, B. (2006). Learning Factors Analysis: A general method for cognitive model evaluation and improvement. *8th International Conference on Intelligent Tutoring Systems*.

Koedinger, McLaughlin, & Stamper (2012). Automated student model improvement. In *Proceedings of the Fifth International Conference on Educational Data Mining*. [Conference best paper.]

Koedinger, Stamper, McLaughlin, & Nixon. (2013). Using data-driven discovery of better student models to improve student learning. *Proceedings of Artificial Intelligence in Education*.

# Learning is complex: Variations in task domains, knowledge demands, student characteristics

- Learning curves showing a decrease in error rate (y-axis) for each successive opportunity (x-axis) to learn
- Averaged across students for different skills – MORE variable



- Averaged across skills for different students – LESS variable



- What causes these variations?

# Turning Discovery into Better Learning

Design
Discover
Deploy
Data



All Selected Knowledge Components

circle-area

compose-by-addition

High rough curve
=> hidden skill
=> redesign instruction
=> Experiment
*Better student learning!*

**Better & faster mastery of solution planning skills**



Instructional Time (minutes)

■ Composition steps
■ Area and other steps

Control: Original tutor

Treatment: Model-based redesign

Koedinger, Stamper, McLaughlin, & Nixon. (2013). Using data-driven discovery of better student models to improve student learning. *Proceedings of Artificial Intelligence in Education.*

# LearnSphere: Integrate across data repositories toward answering questions



**Data collection & analysis time scales**

We need a education data infrastructures to integrate analytic methods
=> *produce discoveries not possible within current data silos*

# Data Integration Example:

*MOOC + OLI = Insight*

What student choices associate with most learning?



*Watching* lecture video

*Reading* web pages

*Doing* online activities with hints & feedback

Learning by doing > 6x better than learning by watching!!

Koedinger et al. (2015). Learning is Not a Spectator Sport: Doing is Better than Watching for Learning from a MOOC. *Proceedings of Learning at Scale*.

# Primary Suggestion for Action: Do data intensive research at our own universities

- Get college instructors involved!
  - Design course activities to collect data
  - Share data & seek analysis partners
  - Engage in discipline-based ed research
- Demonstrate success
  - Set a model for K12
- Incentives
  - NSF fund college-level data-driven innovation
  - Researchers enforce data reuse citation

# Thank you!





Thanks to >200 researchers that have contributed!!

http://learnlab.org/DataShop

Ken Koedinger
koedinger@cmu.edu

# Extras

# Cognitive Model Discovery:
## From qualitative to quantitative

### Traditional Cognitive Task Analysis

- Interviews or think alouds of experts & students
- Result: *Cognitive Model* of expert/student thinking
  - Experts aware of only ~30% of what they know
- Greatly improves instruction
  (~1.5 effect size, Clark et al)



### Data-driven Cognitive Task Analysis

- Use student data from initial tutor
- Goal: more reliable & cost effective
- Employ machine learning & statistics to discover better cognitive models

# Use data to develop models of learners – *because intuition is faulty!*

Which is harder for algebra students?

*Story Problem*

As a waiter, Ted gets $6 per hour. One night he made $66 in tips and earned a total of $81.90.  How many hours did Ted work?

*Word Problem*

Starting with some number, if I multiply it by 6 and then add 66, I get 81.90.  What number did I start with?

*Equation*

x * 6 + 66 = 81.90

Math educators say: story or word is hardest

Students: equations are hardest

**High School Algebra Students**



*Expert blind spot!*
Algebra teachers, especially, incorrectly think equations are easy

Koedinger & Nathan (2004). The real story behind story problems: Effects of representations on quantitative reasoning. *Learning Science.*

# Cognitive Task Analysis using DataShop's learning curve tools



**Without decomposition, using just a single "Geometry" KC,**

**no smooth learning curve.**

**But with decomposition, 12 KCs for area concepts,**

**a smoother learning curve.**

**Upshot**: Can automate analysis & produce better student models

# Discovering a new knowledge component

- Each KC should have:
  - smooth learning curve
  - statistical evidence of learning
  - even error rates across tasks
- Find a feature common to hard tasks but missing in easy ones



1. Not smooth

2. No learning

| KC Name | Intercept | Slope |
|---|---|---|
| circle-area | 0.58 | 0.068 |
| compose-by-addition | 0.74 | 0 |
| compose-by-mult | 0.6 | 0.114 |
| pentagon-area | 0.37 | 0.110 |
| trapezoid-area | 0.35 | 0.091 |

3. Uneven error rate

Easy tasks do not require subgoals, hard tasks do!

# Geometry Tutor
## Scaffolding problem decomposition



Problem decomposition support

# New model discovery: Split "compose" into 3 skills



- Hidden planning knowledge:
  **If** you need to find the area of an irregular shape, **then** try to find the areas of regular shapes that make it up

- Redesign instruction in tutor
  - Design tasks that isolate the hidden planning skill
  - Given square & circle area, find leftover

**3**

When prompts are initially present for component areas

# 3-way split in new model (green) better fits variability in error rates than original (blue)

# Where to go from here?

Possible partnerships/collaborations/relationships to pursue Cyberlearning advances through data sharing?

- Analyses that span levels of analysis

Key needs to be both effective & legal

- Data sharing cyberinfrastructure
    - Easy to use
    - Layered & managed access
    - Rigorous privacy review: IRB+
- Researcher incentives for sharing
    - Sticks: Funder requirements, journal requirements
    - Carrots: Data citation, badges, shared data/analytics counts toward tenure

# What's needed in Cyberlearning data partnerships?

As many as possible of:

- Shared datasets with
  - long-term robust learning & life outcomes
  - multiple assessments: performance, standardized, future learning
  - fine-grain, wide, & deep *click* data
  - fine-grain, wide, & deep *verbal* data
  - embedded experiments: 1 or more random variations
- Analytics sharing with *easy to*
  - access existing analytics
  - apply analytics to full space of Cyberlearning data sources
    - Online courses, simulations, games, tutors, inquiry, class video, ubiquitous computing…
  - recombine existing analytics without programming
  - contribute new analytics & new workflows
- Teams with compatible goals
  - interdisciplinary: education, computer science, psychology, economics …
  - instructors drive research goals
- OTHERS???

# Big Data for Learning Conclusions

- Big data can help unlock mysteries of human learning
  - Science & technology to support learning will transition from Model T to Jet Airplane
- Not the "big" that is important
  - Natural collection: tall, wide, fine, long, deep
- Future: Big data partnerships to tackle big interdisciplinary education questions

# Five Recommendations

1. Search in the "function space"
2. Experimental tests of instructional function decomposability
3. Massive online multifactor studies
4. Learning data infrastructure
5. School-researcher partnerships

Koedinger, Booth, Klahr (2013). Instructional Complexity and the Science to Constrain It. *Science.*