

CS 293/EDUC 473

# Measurement

# Announcements & reminders

- Your reading commentaries are so deep and engaged! Keep up the great work!
- Practice pitches next Monday
  - we'll share your pitch location on Ed
- HW2 due next Tuesday
  - start early especially if you feel uncomfortable with running machine learning models!



# Today's class

- Measurement intro
- Class activity on identifying focusing questions
- Paper discussion on Lee et al. (2024) led by Nic, Cameron and TJ

# Where are we?



Identify  
Problem



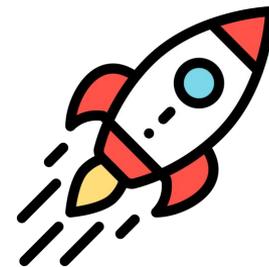
Data  
Exploration



Algorithm  
Development  
& Validation



Tool  
Development



Deployment

Overarching Themes:



Bias &  
Fairness



Working  
closely with  
teachers

## What would you like more dedicated class discussion on?

Nobody has responded yet.

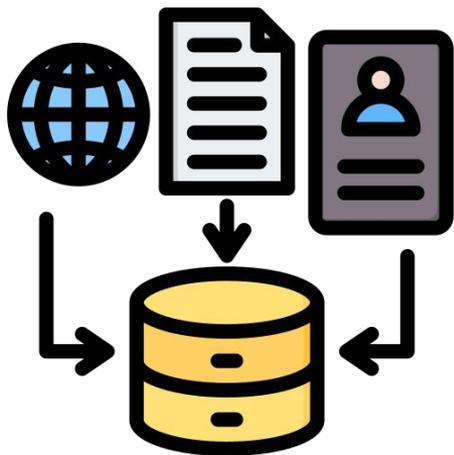
Hang tight! Responses are coming in.



# How do we define measurement?

## Data

Structured (e.g. likert scale responses) /  
Unstructured (e.g. language)



## Score / Label

Measuring a target construct

**Binary:** Is this utterance on task?

**Continuous:** To what extent does the student feel empowered in this classroom?

**Categorical:** What is the topic of this lesson?

You can use these in  
quantitative analyses:



**Most aspects** of a quantitative research project / intervention / tool require measurement

- Identifying & analyzing teaching practices
- Evaluating fairness & bias
- Identifying need for intervention
- Understanding teachers' and students' perceptions of a tool
- Measuring outcomes
- ...

**None of these are trivial to measure**

# Most aspects of a quantitative research project / intervention / tool require measurement

## Ex. Challenges and Issues

- Identifying & analyzing teaching practices
  - Evaluating fairness & bias
  - Identifying need for intervention
  - Understanding teachers' and students' perceptions of a tool
  - Measuring outcomes (e.g. student learning)
  - ...
- Subjectivity and context-dependence
- Sparse data & oversimplification of demographic categories
- High stakes & prone to bias
- Low response rate & self-reporting bias
- Choice of outcomes are often the most controversial

## What type of NLP measures does your final project require?

Nobody has responded yet.

Hang tight! Responses are coming in.



# Do you need to identify the measurement target?

E.g., type of classroom practice, dimension for user attitude (e.g. difficulty).

No:

Skip to next step

Yes:

Pick the target at the intersection of **promise & feasibility**

- Lit review
- Talk to people
- Look at data



## Example brainstorming spreadsheet (list of discourse practices relevant to math ed)

Category	Feature Associated with Better Learning Outcomes	What to Measure	Examples
General pedagogical	student participation	number of words uttered by students/minute; and/or length of student utterances	
General pedagogical	test-orientedness	measure the number of references to standardized testing	MCAS, DCCAS
General pedagogical	instructional time	amount of instructional time vs off task time	procedural talk vs instructional talk; noise in the classroom
General pedagogical	wait time after questions	amount of wait time after questions	
General pedagogical	check-in on students	number of times teacher asks questions that check in to see if students are following along	"make sense?" "Any questions"? Thumbs up? Hold your boards up; "student .. looks puzzled"
Math specific	use of math terms	density of math terms from teachers & students; measure to what extent teachers press students to use such terms	angle, fraction
Math specific	use of sloppy (math) terms	density of sloppy terms from teachers and from students	borrowing, top and bottom, cancelling
Math specific	use of proofs, mathematical reasoning/explanation	presence of proofs, mathematical reasoning/explanation in teacher & student talk	
Math specific	degree of direct instruction & focus on memorization	estimate the degree to which the teacher is doing direct instruction	teacher talk with short student answers interspersed; words like remember, recall; first thing you do when you...what do you do next...we're also going to have to do what?
Math specific	cognitive demand of (math) questions	estimate the degree of cognitive demand of questions (teachers & students)	why, explain, what does that mean, different, difference, compare, what's missing, how do these relate
Math specific	teachers' evaluation of student contributions	see whether and how the teacher remediates students' misunderstandings;	correction, reformulation, repetition, praise
Math specific	uptake	degree to which teacher uses students' mathematical contributions in subsequent instruction; students' uptake of other students' ideas (with minimal teacher orchestration)	
Classroom climate	positive references to students	degree to which teacher uses student names in a positive way	positive: "Geoffrey's idea" or "Marie, tell us what you are thinking" "I think Nonie solved the problem in the same way"; negative: Geoffrey!
Classroom climate	affirmation of knowledge and skill	degree to which teacher encourages students	"You totally understand this, you just need to tweak what you're saying a little bit"
Classroom climate	broad regard	degree to which the teacher shows interest in the students' lives	asking non-academic questions

# Does an NLP measure exist already for what you want to do?

Yes:

Skip to validation (on your domain)

No:

**Develop a measure** (in most cases) following the standard paradigm → next slide

# Standard NLP measure development workflow

## 1. Create high quality validation set

- a. With sufficient # of examples to capture relevant variation (rule of thumb: at least 1k examples for a relatively straightforward measure, 2k for more subjective ones)
- b. When possible, create a held-out test set too (that you only evaluate on at the very end)

## 2. Iteratively develop & validate model

- a. Supervised paradigm: label training data → train classification/regression model
- b. Unsupervised/self-supervised paradigm: leverage unlabeled data

# Standard NLP measure development workflow

## 1. Create high quality validation set

- a. With sufficient # of examples to capture relevant variation (rule of thumb: at least 1k examples for a relatively straightforward measure, 2k for more subjective ones)
- b. When possible, create a held-out test set too (that you only evaluate on at the very end)

**what if creating a validation set is not at all trivial because the construct is highly subjective?**

# What to do if your **interrater agreement is fair to moderate?**

Even when working with domain experts & doing several rounds of rater training and discussion

## First ask: why is agreement low?

Potential cause	Potential solutions
Poorly defined construct	<ul style="list-style-type: none"><li>● Improve definition &amp; coding scheme!</li><li>● Revise input level of granularity (e.g. sentence vs utterance vs segment)</li></ul>
Context-dependence of construct	<ul style="list-style-type: none"><li>● (When possible) Add more context</li><li>● (When appropriate) Pre-define context</li></ul>
Intersubjectivity (diff. people might perceive or react to the same thing differently)	This may be important variation that's worth preserving

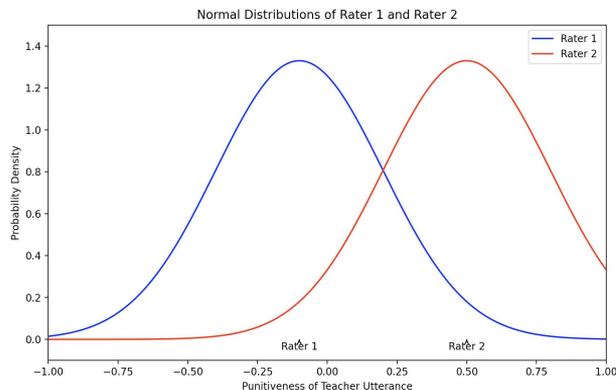
# How to handle **inherently subjective** constructs?

During annotation

- Have multiple annotators (the more the better) for each example.

Processing Annotations

- Z-score judgments before aggregating
- Create different subsets of the data (subjective subset; less subjective subset) for evaluation



# How to handle **inherently subjective** constructs?

## During annotation

- Have multiple annotators (the more the better) for each example.

## Processing Annotations

- Z-score judgments before aggregating
- Create different subsets of the data (subjective subset; less subjective subset) for evaluation

## Modeling

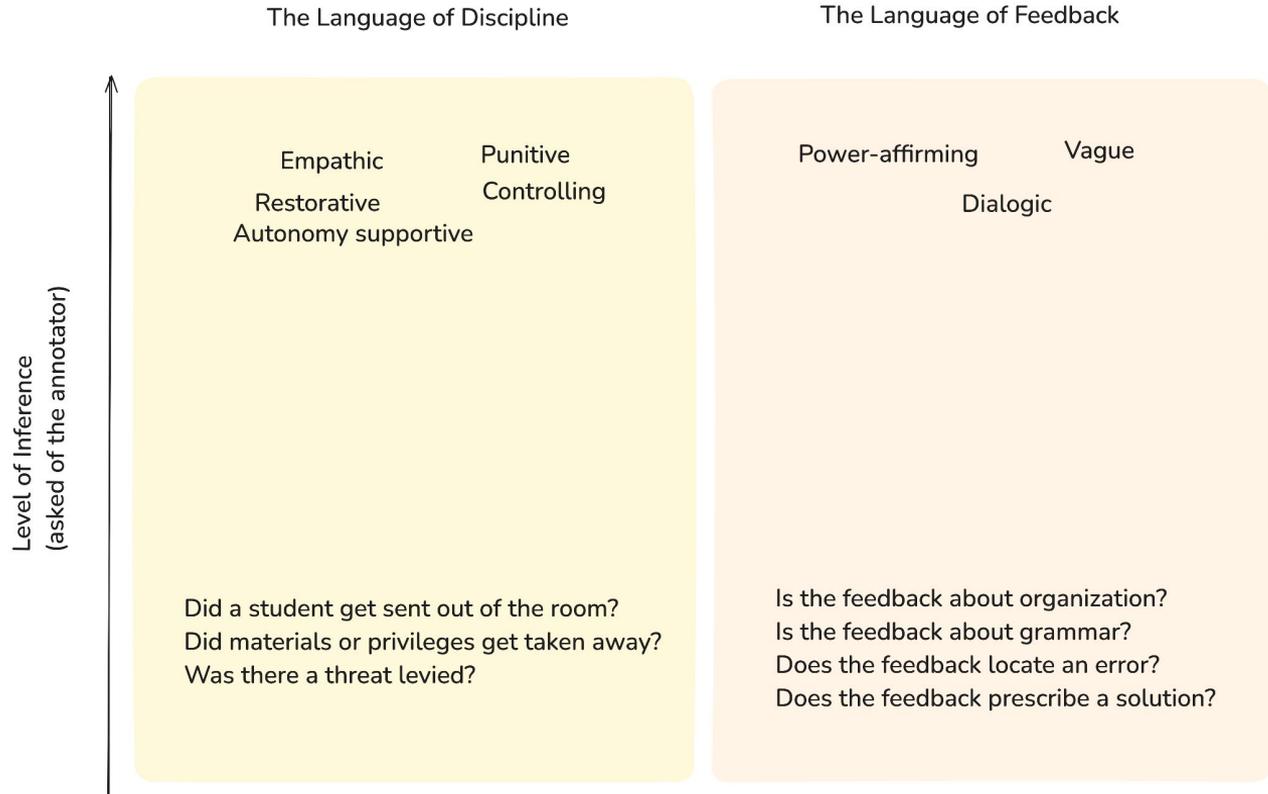
- Incorporate confidence into the measurement
  - e.g. build a model that predicts rater agreement as a proxy for confidence
- Train representations separately for each rater and then combine them into a shared representation ([Davani et al., 2022](#))

## Validation

- Check if results are robust to variations in data or modeling decisions
  - e.g. leave-out validation
- Don't rely too heavily on your "ground truth" values
  - Correlate your measure with other relevant variables (e.g. overall instruction quality) to understand if it relates to positive or negative outcomes
  - Estimate the impact of applying your measure to address a specific issue
- Don't use for making high stakes decisions!

## Application

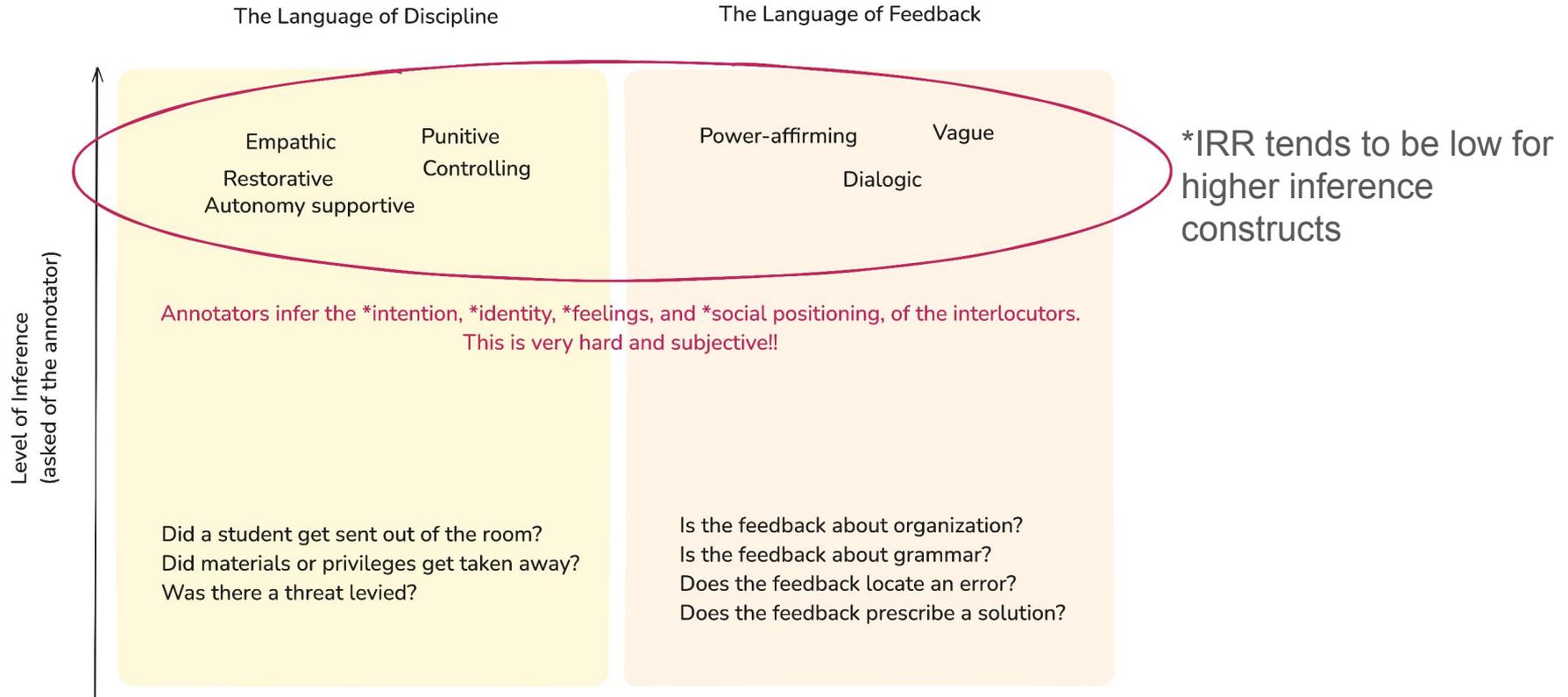
# A few examples from lab work:



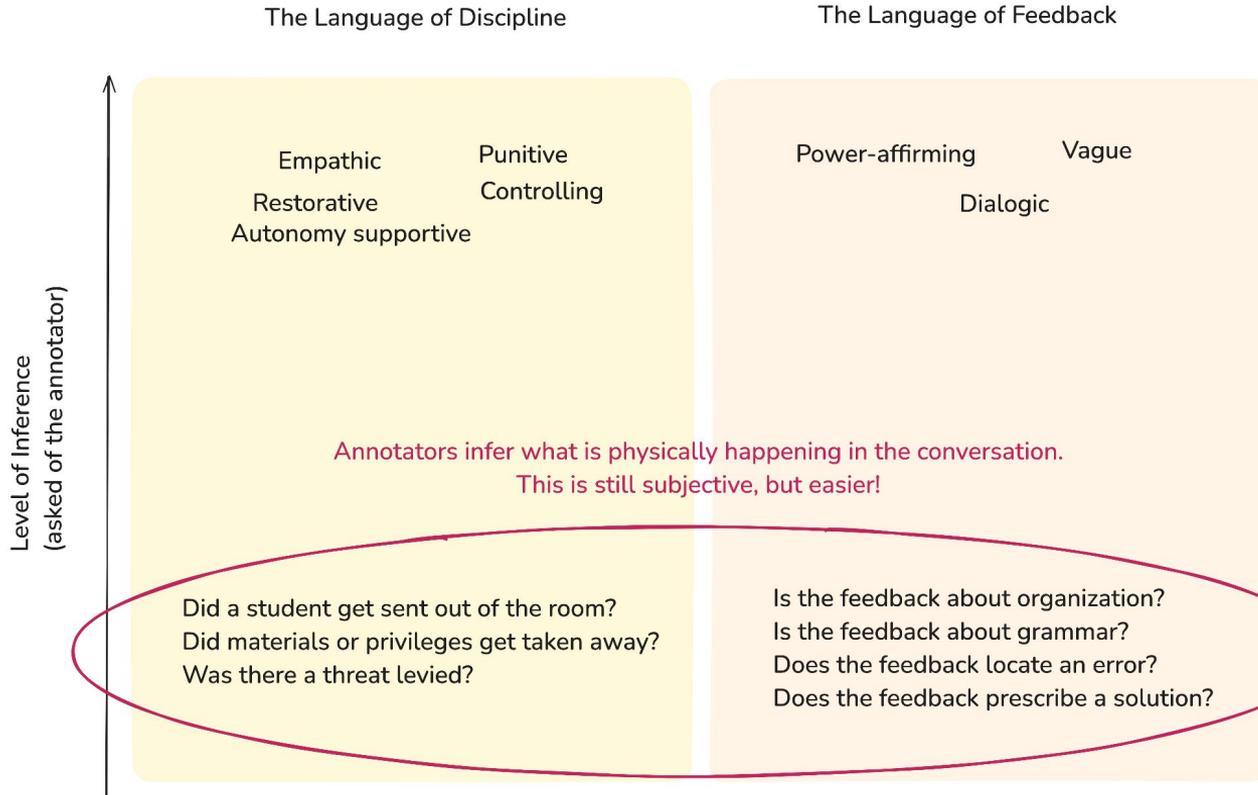
\*all of these words are from the literature!

domain problems != technical problems

# A few examples from lab work:



# A few examples from lab work:



*Do we expect LLMs to have the context to make this judgement based on its training? Often for conversational data... we think no.*

\*is the construct so literal that we trust an LLM to infer instead?

# Supervised modeling: LLMs or smaller models?

Smaller models (RoBERTa, BERT, etc.)	LLMs
Resources: <a href="https://simpletransformers.ai/">https://simpletransformers.ai/</a> ; <a href="https://huggingface.co/docs/transformers/index">https://huggingface.co/docs/transformers/index</a>	GPT, Claude, Gemini, DeepSeek
<p>Pros:</p> <ul style="list-style-type: none"><li>● Downloadable → more transparency &amp; control</li><li>● Needs little compute; <b>greener</b> choice!</li><li>● Can achieve similar performance to LLMs when sufficient labeled data is available</li></ul>	<p>Pros:</p> <ul style="list-style-type: none"><li>● Very good at few shot learning</li><li>● Can be tuned with instructions</li></ul>
<p>Cons:</p> <ul style="list-style-type: none"><li>● Require more training data</li><li>● Can't be tuned with instructions or via interacting with the model</li></ul>	<p>Cons:</p> <ul style="list-style-type: none"><li>● Most cannot be downloaded</li><li>● Some models can't be finetuned</li></ul>

Superv

or both?

Smaller

**model distillation**: use LLMs to generate data for finetuning smaller models

Resource

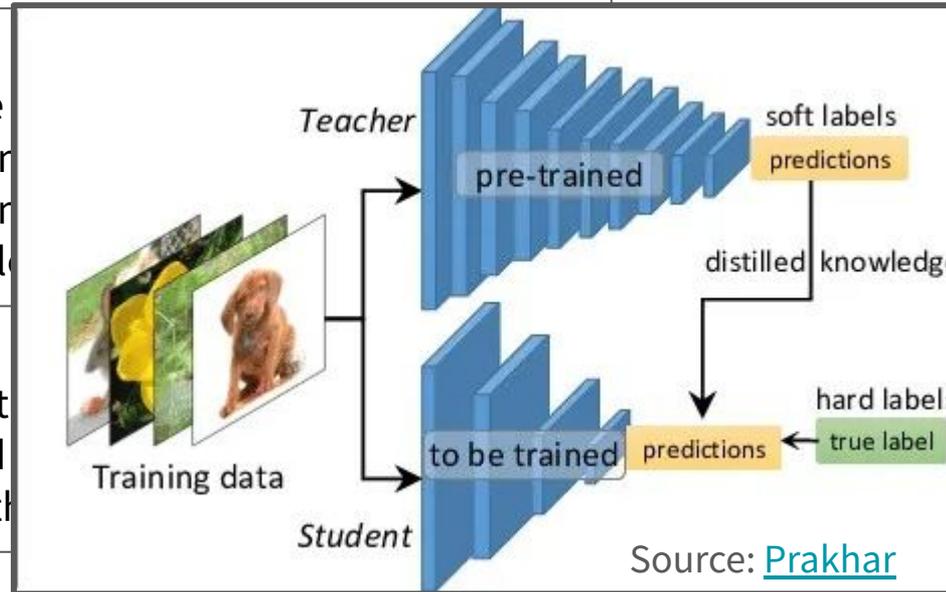
<https://h>

Pros:

- Downloadable
  - Needs little con
  - Can achieve sim
- sufficient labels

Cons:

- Require more t
  - Can't be tuned
- interacting with



few shot learning with instructions

be downloaded  
s can't be finetuned

Source: [Prakhar](#)



# Edu-ConvoKit

## A Pre-Processing

## B Annotation

## C Analysis

### Dataset

Speaker	Text
Tutor	Alice, what units would we use to measure speed?
Student	Miles per hour.
...	...

e.g., tutoring or classroom conversation

Speaker	Text
Tutor	[STUDENT] what units would we use to measure speed?
Student	Miles per hour.
...	...

e.g., anonymization, grouping utterances

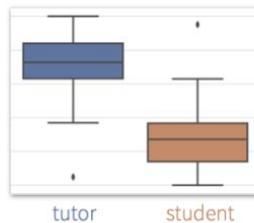
Speaker	Text	Feature
Tutor	[STUDENT], what units would we use to measure speed?	
Student	Miles per hour.	
...	...	

e.g., talk time, student reasoning, teacher talk moves

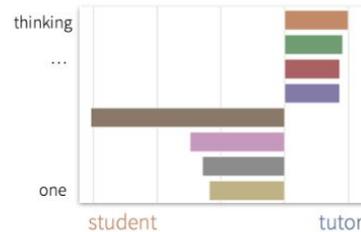
**Qualitative** e.g., student reasoning examples

Student: Because, um, because, let's say at this point you're not traveling at fifteen miles an hour [...]

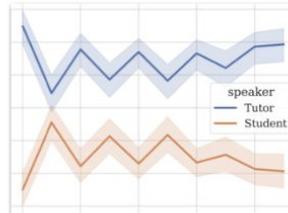
**Quantitative** e.g., talk time



**Lexical** e.g., log odds



**Temporal** e.g., talk over time



**GPT** e.g., summary of conversation

🤖: The tutor asked the student to determine which is larger, one half or two thirds, and by how much. The students used different colored rods to represent fractions [...]

# Class activity: identifying **focusing questions**

Focusing questions press students to communicate their thoughts and reflect on those of their classmates.

Teacher:  $(0,0)$  and  $(4,1)$  are two points on a line. What's the slope?

(possible follow up questions)

funnel

Teacher: What's the rise? What's the run?



Students: 1, 4

focus

Teacher: What do you think of when I say slope?



Student: The angle of the line.

Student: Fractions.

Student: How fast the line changes.

## Examples of Common Types of Focusing Questions

- **What** - “What else could you try?” “What do you think about that?” “What does your answer mean?” “What else do we know about \_\_\_?”
- **How** - “How did you come up with that solution?” “How else might you think about the problem?” “How did you know where to start?” “Did you pull that [number/operation] from the [representation/word problem]?”
- **Why** - “Why do you think that method works?” “Why did you start with that part?” “Is that estimate reasonable?”
- **Tell me more** - “Okay. Go ahead. Tell me more about the strategy you chose.”
- **Comparing to similar/different methods** - “How is this similar/different to Student A’s method?” “That’s an interesting method. Does it work here?”
- **Comparing to similar concepts** - “How are the attributes of [geometric shape] similar to the [other geometric shape] we learned last week?”
- **Possibility of New Scenarios** - “What might happen if \_\_\_?” “What else would have to happen?”
- **Responding to Peers** - “I hear her saying \_\_\_\_. Does anyone else know what she means when she says that?” “Can someone revoice what Student Z just said?”

Reading discussion on Lee et al. (2024)