

Quality Control - part 2

Crowdsourcing and Human Computation

Instructor: Chris Callison-Burch

Website: crowdsourcing-class.org

Different Mechanisms for Quality Control

- Aggregation and redundancy
- Embedded gold standard data
- **Economic incentives**
- Reputation systems
- Statistical models

Does pay impact quality?

- Economic theory holds that workers are rational actors
- Will choose to improve their performance in response to a scheme that rewards improvements with financial gain
- Example: executive compensation tied to stock price

Different pay schemes

- Lazear studied of workers who installed windshields on a production line
- Switched from pay per hour to pay per unit during a year and a half
- Individual productivity for workers who started in the hourly rate and switched to the per-unit scheme increased by 20%
- Conclusion: performance-based pay schemes can elicit improved performance

Is that the whole story?

- Sometimes financial incentives can undermine “intrinsic motivation”. This can lead to poorer outcomes.
- For complex tasks, performance pay can encourage workers to focus only on the aspects of their jobs that are actively measured
- Can also lead to employees avoid taking risks, thereby hampering innovation

Financial Incentives and the “Performance of Crowds”

- Experiment with economic incentives on Amazon Mechanical Turk
- An exciting tool for behavioral research, since you can recruit thousands of participants from a real labor market

Impact of compensation

- Does compensation change the quantity of work performed (output)?
- Does it change the quality of the work (accuracy)?

Re-order Traffic Images

Unsorted



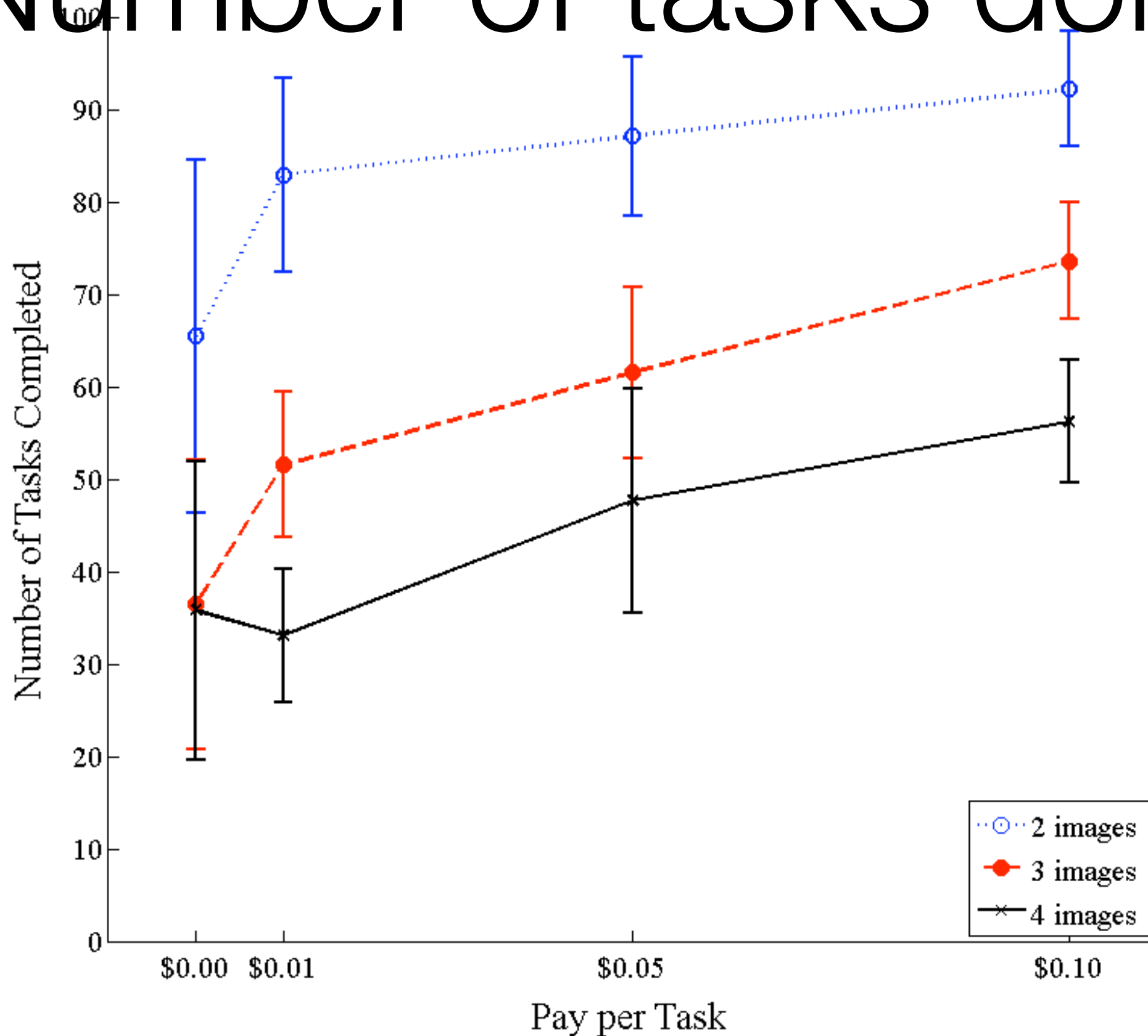
Sorted



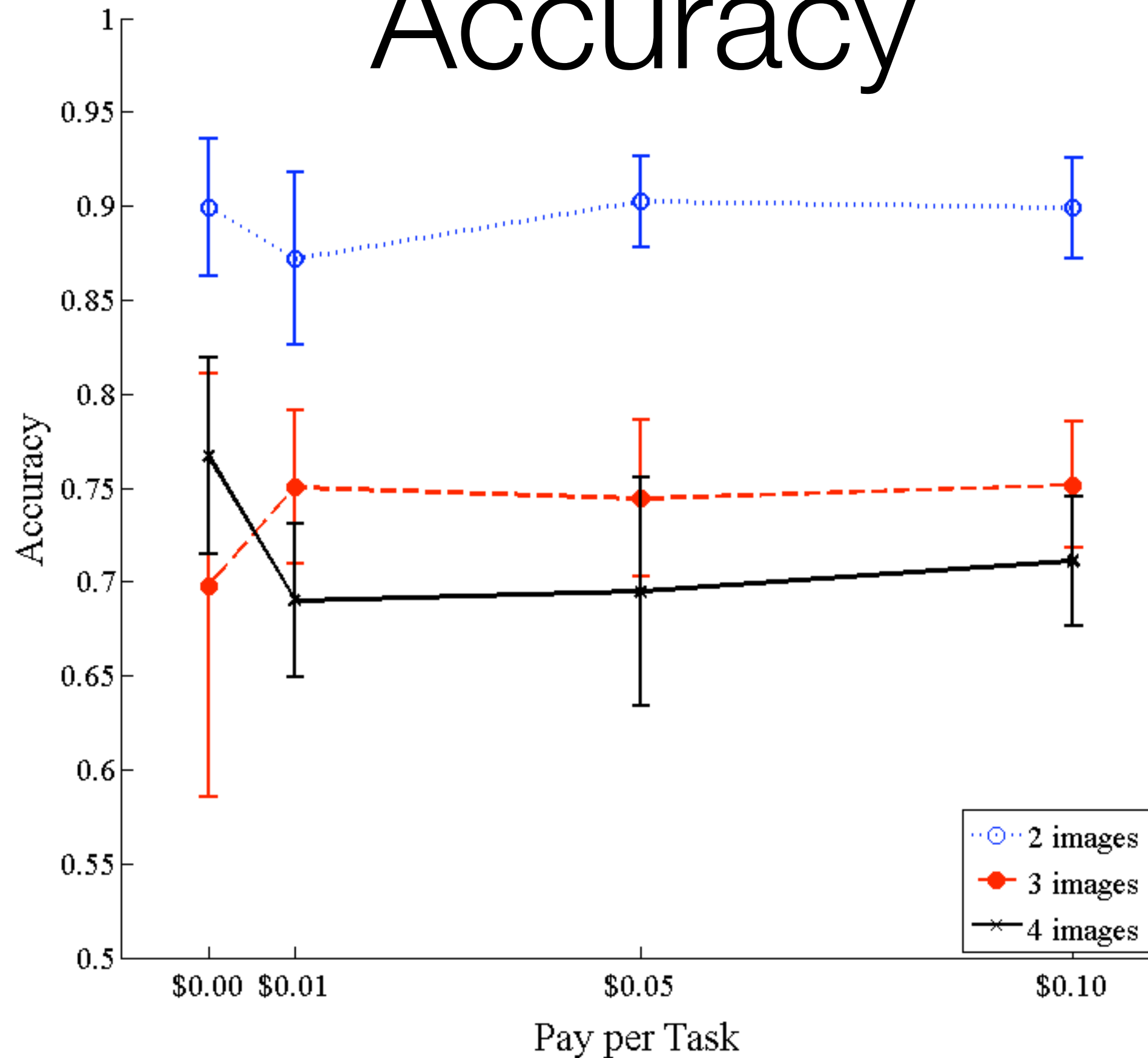
Payment scheme

- Everyone: \$0.10 for doing training examples and filling out a survey
- Payment levels: nothing, 1¢, 5¢, 10¢ per set
- Num images per set (independent of payment): 2, 3, 4
- Each person sorted up to 99 sets of images, could end participation at any point and get paid for what they did
- 611 subjects sorted a total of 36,425 image sets

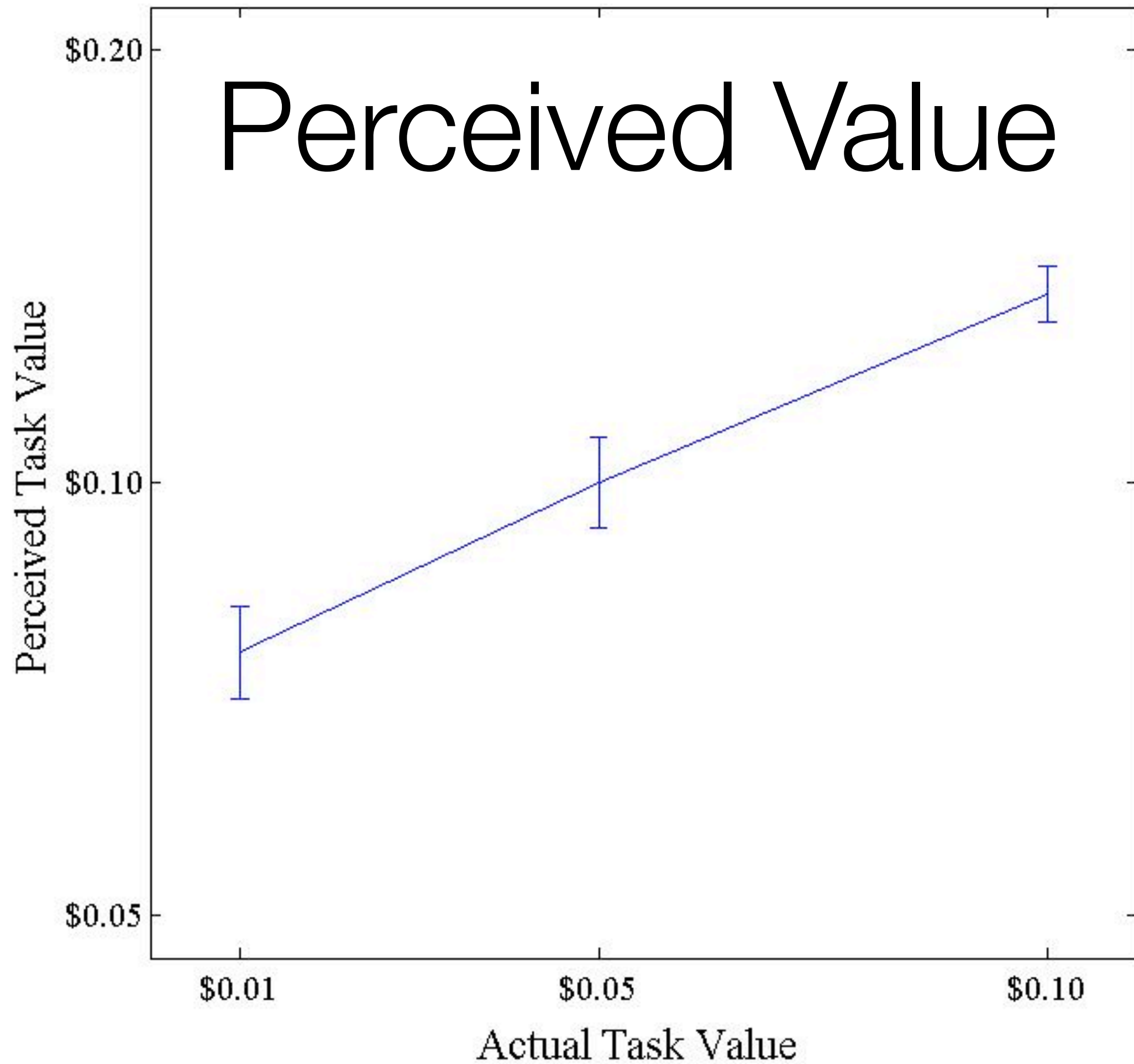
Number of tasks done



Accuracy



Perceived Value



Word Jumble Puzzles

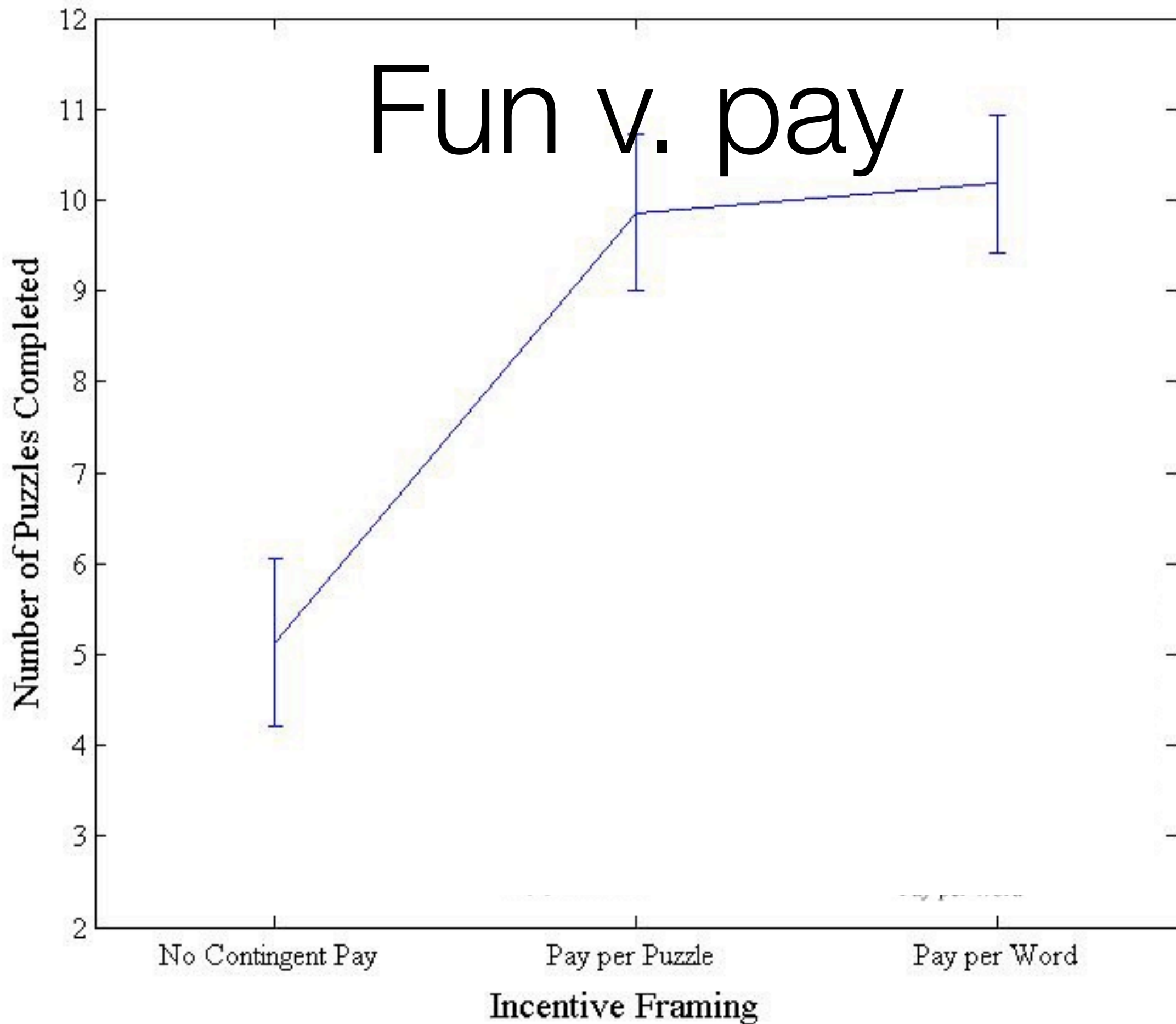


- Find as many of the of words in a set as you can:
- ACHIEVE, ATTAIN, BUILDING, CHAIR, COMPLETE, GREEN, LAMP, MASTER, MUSIC, PLANT, STAPLE, STEREO, STRIVE, SUCCEED, TURTLE
- Not all of the words listed are in the puzzle!

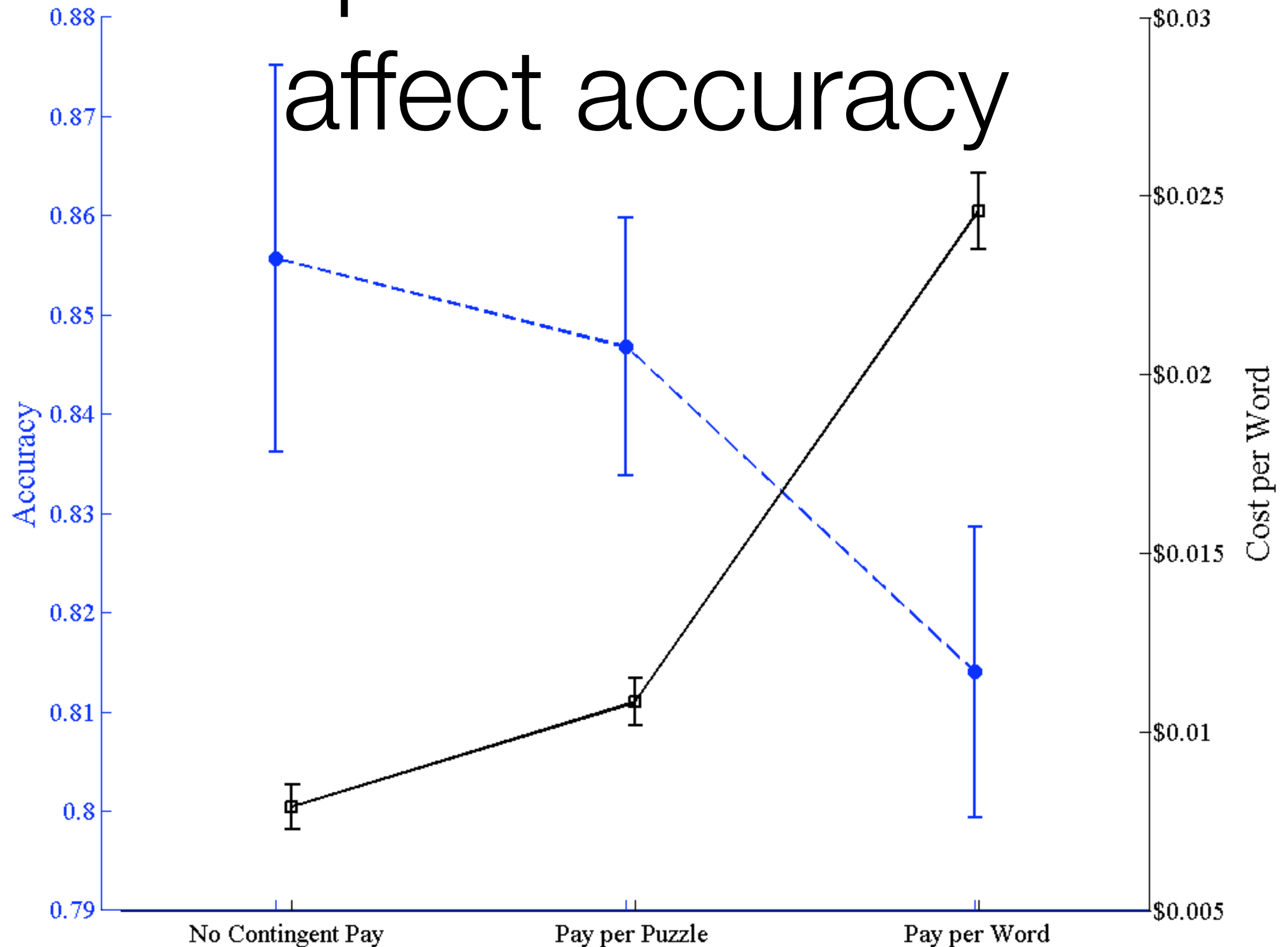
Experimental setup

- Different pay rates (just as before)
- Subjects were told that they would be paid either on a per-grid basis or a per-word basis, or not told anything
- quantity = number of puzzles completed
quality = fraction of words found per puzzle
- Participants could do up to 24 puzzles
- 320 subjects solved 2736 puzzles, finding 23,440 words

Fun v. pay



Compensation doesn't affect accuracy



Perceived Value

\$0.10

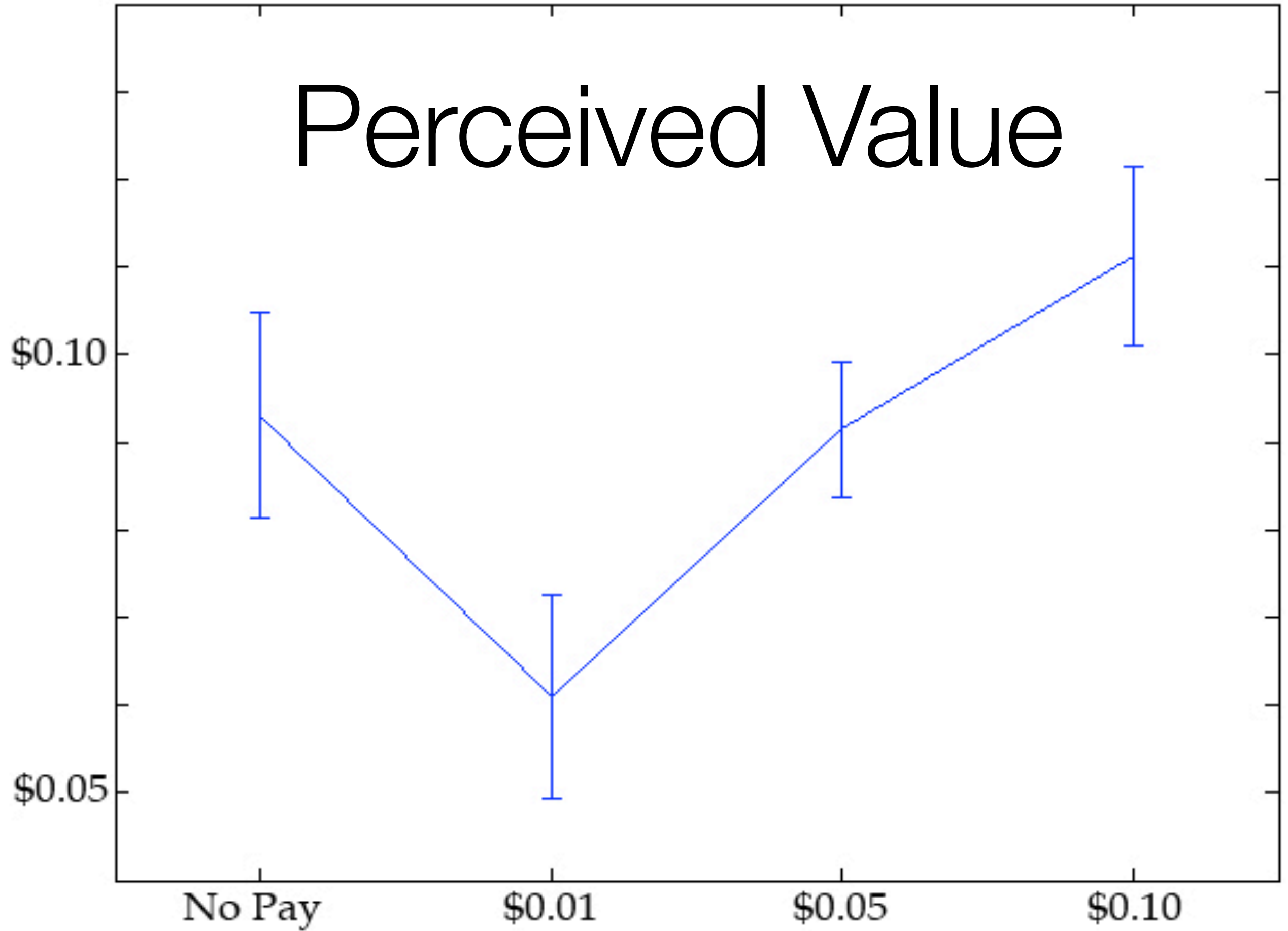
\$0.05

No Pay

\$0.01

\$0.05

\$0.10



Findings

- Paying subjects elicited higher output than gamification, and increasing pay rate yielded even higher output
- However, paying subjects did not affect their accuracy
- Anchoring effects are significant – the reward you set impacts perceived value

Implications for your tasks?

- When you can use non-financial rewards, like intrinsic motivation, do so, since the quality of work will be the same
- When you can't use intrinsic motivation, it might be in your best interest to pay as little as possible. Your work will be done slower, but quality will be similar.
- Is this fair to workers?

What do you think?

- Is studying workers on Mechanical Turk a valid way of studying other labor markets?
- What possible confounds are there?
- What could we do to control for them?