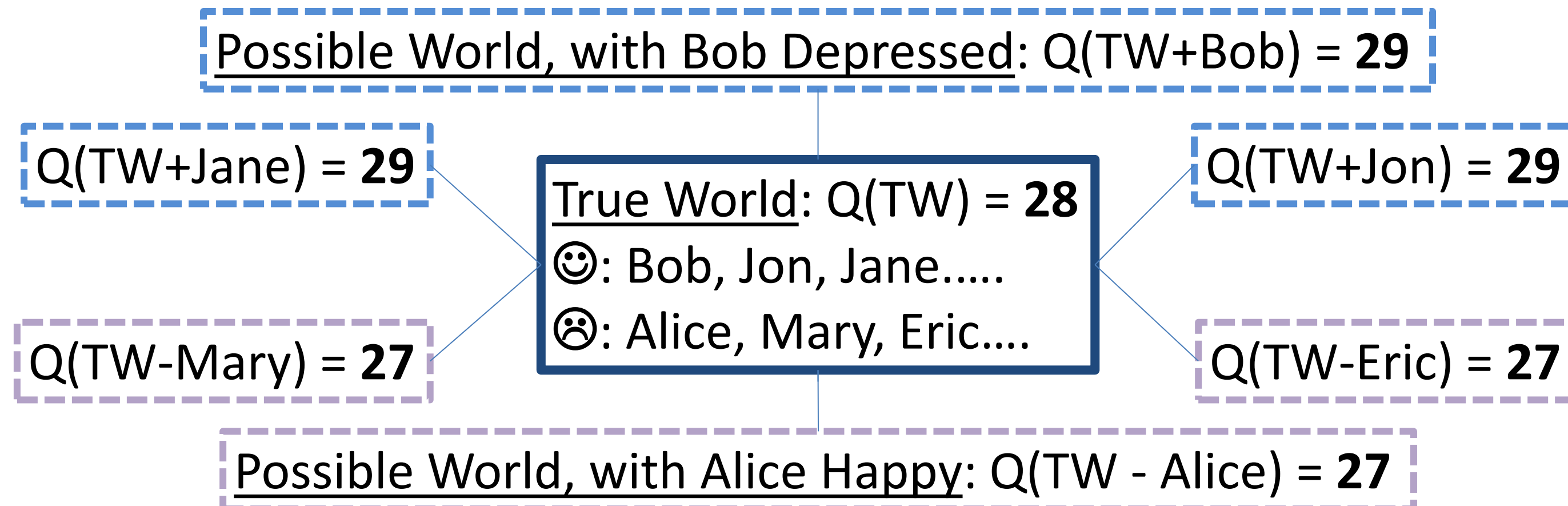


Publicly Constrained Populations in Differential Privacy

Christine Task, Chris Clifton, Computer Science Department, Purdue University

Differentially private queries use noise to blur borders of neighboring possible worlds:

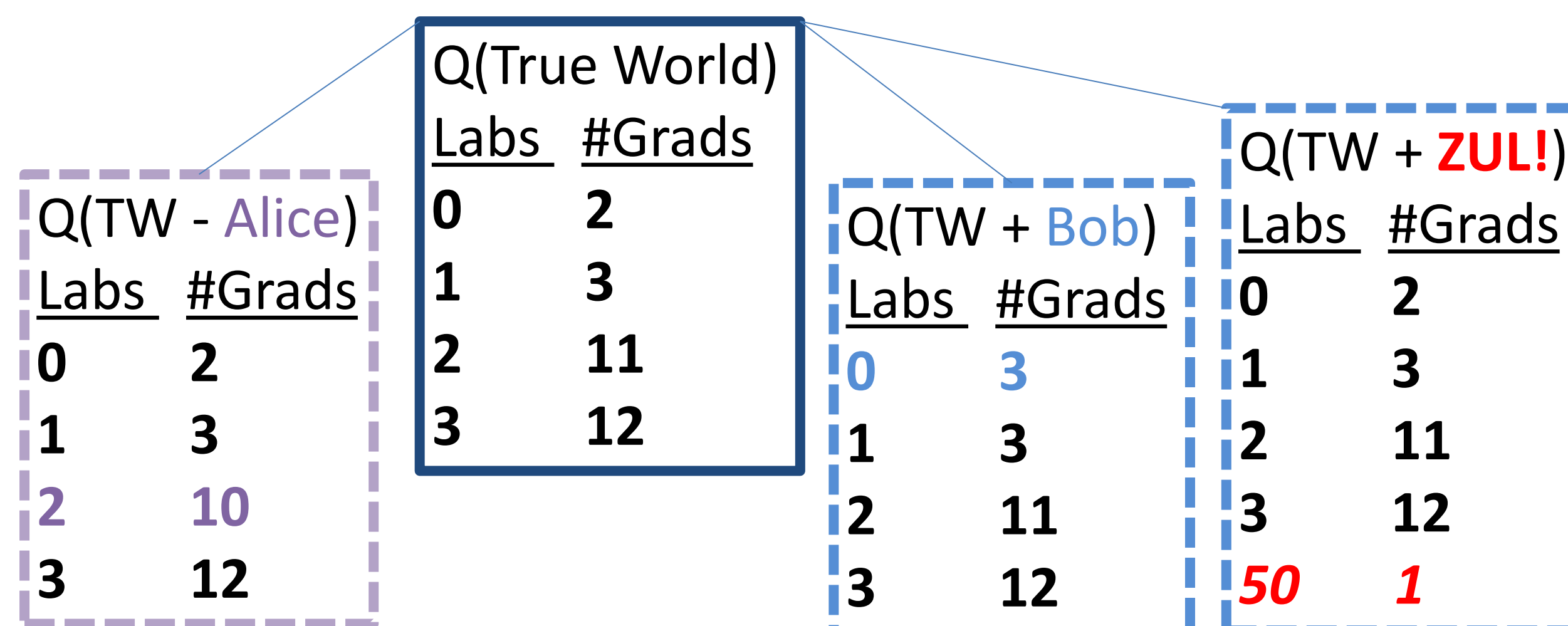
Q: "How many Purdue CS graduate students report seeking treatment for depression?"



$Q_{\text{private}}(\text{TW}) = 27.2$
 (Add random noise of standard deviation 1 to Q)
 How are Alice, Bob? We can't tell: The result 27.2 could have come from any neighbor world.

Which pairs of worlds do we consider to be neighbors? Traditionally: Every pair of data-sets that differ in a single *arbitrary* individual data-point.

Query: "What is the teaching load of Purdue CS graduate students who report depression?" (results as histogram over # of labs taught)



Labs	#Grads
0	.9
1	5.1
2	10.7
3	13.3
...	
50	.8
51	.2
52	1.1
...	

What if it's public knowledge that some data-points simply *can't occur* in the data-set?

What does an attacker learn from the privatized query result R ? For any neighbors X, Y :

$$\frac{\text{Prob}(\text{True World } (TW)=X | Q_{\text{private}}=R)}{\text{Prob}(\text{True World } (TW)=Y | Q_{\text{private}}=R)} = L_{x,y} \frac{\text{Prob}(TW=X)}{\text{Prob}(TW=Y)}, \text{ where } L_{x,y} = \frac{\text{Prob}(Q_{\text{private}}=R | TW=X)}{\text{Prob}(Q_{\text{private}}=R | TW=Y)}$$

$L_{x,y}$ measures how much an attacker learns from R about the relative likelihood of X vs. Y : which world is more likely true. Differential Privacy requires that L be fairly close to 1, $L \leq e^\epsilon$ for small $\epsilon > 0$. If $Q(X) - Q(Y)$ is large, then large noise is needed to ensure this.

But if Y has an arbitrary, impossible neighbor Z :

$$\frac{\text{Prob}(TW=Z | Q_{\text{private}}=R)}{\text{Prob}(TW=Y | Q_{\text{private}}=R)} = L_{z,y} \frac{0}{\text{Prob}(TW=Y)} = 0$$

We can ignore $L_{z,y}$ when calibrating noise, because the attacker's belief in the relative likelihood of X vs. Z cannot change. She knows $\text{Prob}(Z) = 0$.

Many queries that appear to require too much noise to be feasibly privatizable in general, can be privatized with fairly little noise when public knowledge is considered.

Query	Generic Noise	Potential Context	Constrained Noise
How many triangles are there in graph G ?	Unbounded	Given G is a subset of public school twitter follower graph that represents sexual relations.	Bounded by max individual triangle count of public graph
Average GPA in CS426?	Noise in denominator, orders of magnitude error	Given public knowledge that n people are registered for the class, no noise necessary in denominator.	Standard deviation of noise values is just $1/n$
How many in social network S have more than k friends?	Unbounded	Given that S limits members to at most m friends.	Standard deviation of noise is m , good when $ S \gg m$