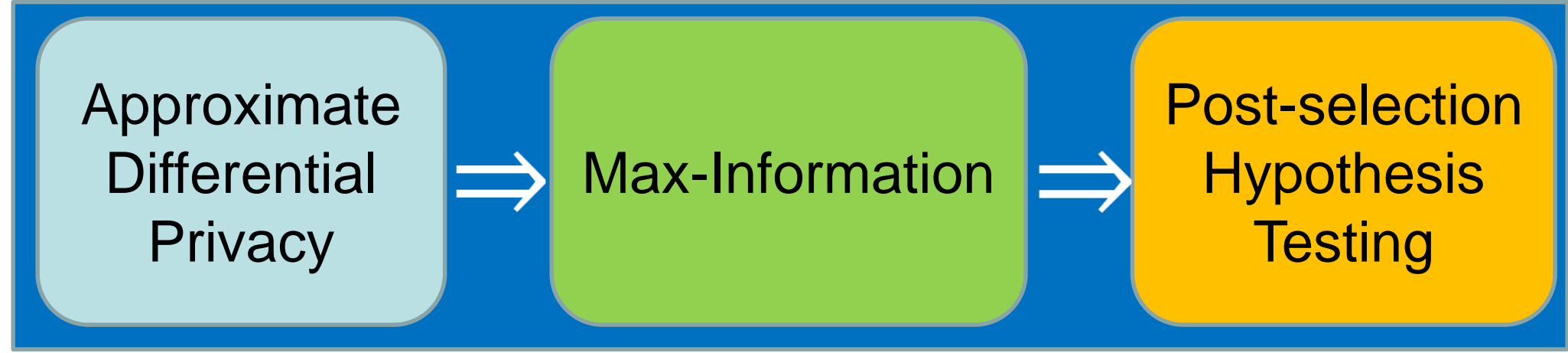
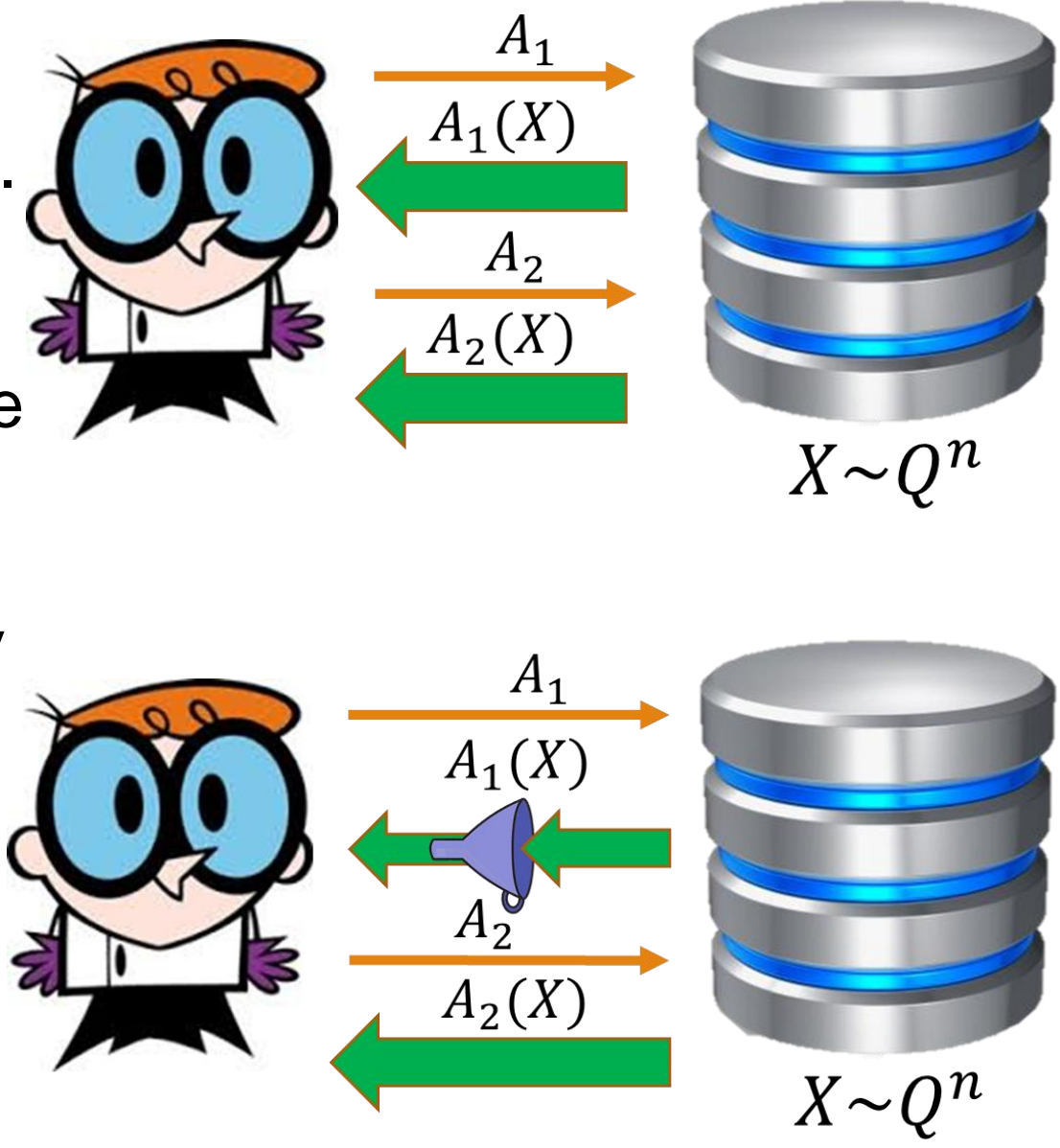


*Adaptive Data Analysis* refers to the reuse of data to perform analyses suggested by the outcomes of previously computed statistics on the same data. In this work, we initiate a principled study of how the generalization properties of *approximate differential privacy* can be used to perform *adaptive hypothesis testing*. This substantially extends the existing connection between differential privacy and *max-information*, which previously was only known to hold for pure differential privacy. It also extends our understanding of max-information as a partially unifying measure controlling the generalization properties of adaptive data analyses.



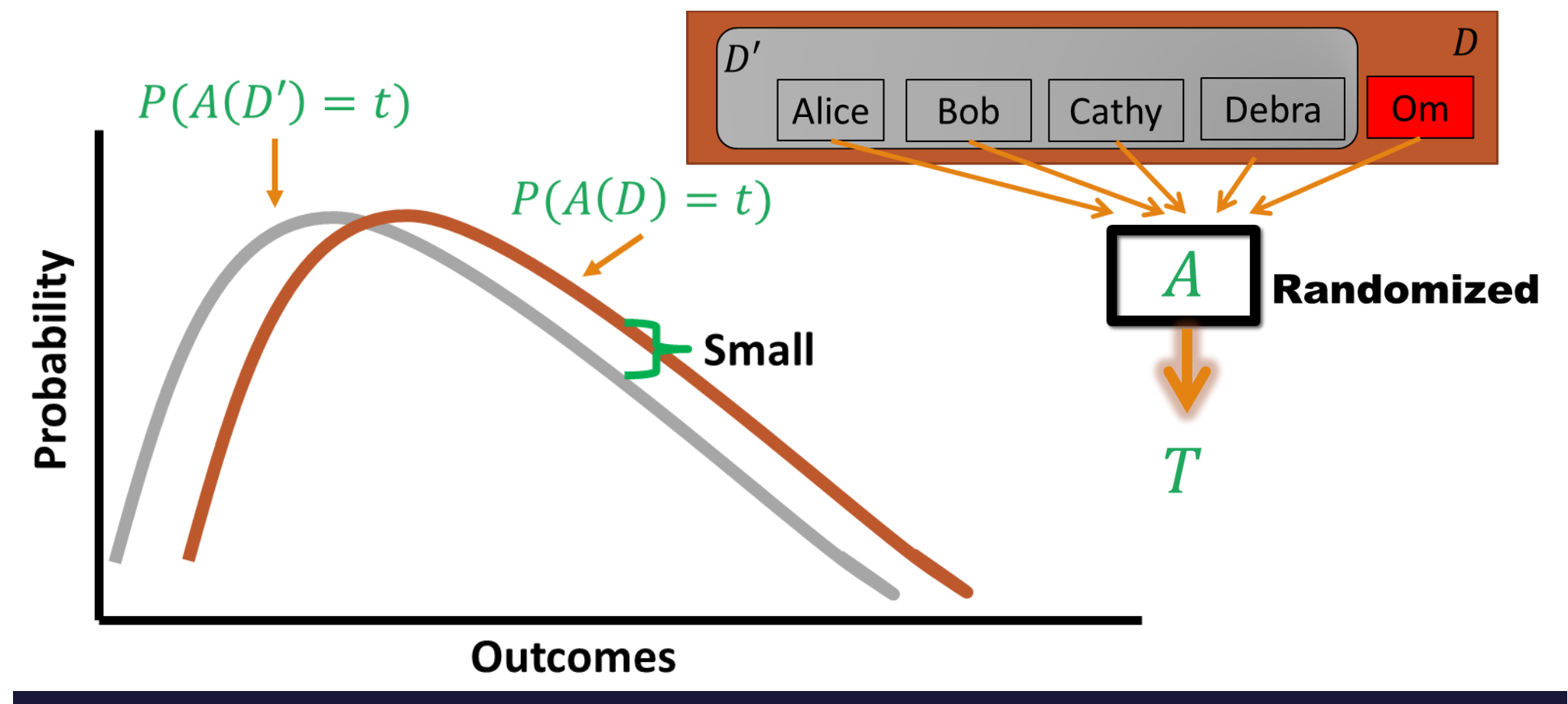
## Adaptive Data Analysis

- A lot of existing theory assumes tests are selected independently of the data.
- In practice, data analysis is inherently interactive, where experiments may depend on previous outcomes from the **same** dataset.
- **Question:** How can we provide statistically valid answers to adaptively chosen analyses?
- **Answer:** **Limit** the information learned about the dataset. [DFH+15a]
- Part of a line of work initiated by [DFH+15a, DFH+15b, HU14].



## Differential Privacy [DMNS06]

- A randomized algorithm  $A: D^n \rightarrow T$  is  $(\epsilon, \delta)$ -differentially private if for all neighboring data sets  $x, y \in D^n$ , i.e.,  $dist(x, y) = 1$ , and for all sets of outcomes  $S \subseteq T$ , we have
 
$$P(A(x) \in S) \leq e^\epsilon P(A(y) \in S) + \delta$$
- If  $\delta=0$ , we say pure DP. If  $\delta>0$ , we say approximate DP.



## Post-Selection Hypothesis Testing

- Hypothesis test: Defined by a null hypothesis  $H_0$  and a test statistic  $t$ .
- Purpose: **Reject**  $H_0$  if the data  $X$  is **not likely** to have been generated from some distribution  $Q^n$  such that  $Q \in H_0$ .
- Significance level of  $t = \alpha \Rightarrow \Pr_{X \sim Q^n}[t(X) = \text{Reject}] \leq \alpha$ .
- Assumes choice of  $t$  is independent of the data  $X$ .
- Goal: For an adaptively chosen test  $t_{A(X)}$ , we want to bound
 
$$\Pr_{X \sim Q^n}[t_{A(X)}(X) = \text{Reject}]$$
 for  $Q \in H_0$ .
- Problem:  $t_{A(X)}$  can be tailored **specifically** to  $X$ .

## Technical Contributions

- Previous results [DFH+15a]: If  $A: D^n \rightarrow T$  is  $(\epsilon, 0)$ -DP,
  - For  $\beta > 0$ , we have  $I_{\infty, \Pi}^\beta(A, n) \leq \tilde{O}(\epsilon^2 n)$
  - $I_\infty^0(A, n) \leq \epsilon n$
- **Positive Result:** If  $A: D^n \rightarrow T$  is  $(\epsilon, \delta)$ -DP, for  $\beta \approx O(n\sqrt{\delta/\epsilon})$ ,
  - we have  $I_{\infty, \Pi}^\beta(A, n) = O(\epsilon^2 n + n\sqrt{\delta/\epsilon})$



- Consequences:
  - $k$  rounds of adaptivity: max-information  $\sim k$  rather than  $k^2$
  - Generalizes and unifies previous work
- **Negative Result:**  $\exists$  an  $(\epsilon, \delta)$ -DP algorithm  $A$  s.t.
 
$$I_\infty^\beta(A, n) \approx n$$
 for any  $\beta \leq \frac{1}{2} - \delta$ .

## Max-Information [DFH+15b]

- An algorithm  $A$  with **bounded max-info** allows the analyst to treat the output  $A(X)$  as if it is **independent** of data  $X$  up to a factor.
 
$$I_\infty^\beta(X; A(X)) \leq k \Leftrightarrow \Pr_{(x,a)} \left( \log \left( \frac{\Pr[X=x, A(X)=a]}{\Pr[X'=x] \Pr[A(X)=a]} \right) > k \right) \leq \beta$$
- Differentiate between general and **product** distributions:
 
$$I_\infty^\beta(A, n) = \sup_{S: X \sim S} I_\infty^\beta(X; A(X))$$

$$I_{\infty, \Pi}^\beta(A, n) = \sup_{P: X \sim P^n} I_\infty^\beta(X; A(X))$$
- [RRST.16]: If  $I_{\infty, \Pi}^\beta(A, n) \leq k$ , then for  $\gamma(\alpha) = \frac{\alpha - \beta}{2^k}$ ,
 
$$\text{Significance level of } t_{A(X)} = \gamma(\alpha) \Rightarrow \Pr_{X \sim Q^n}[t_{A(X)}(X) = \text{Reject}] \leq \alpha.$$



## Related Publications

- [BNS+16] Raef Bassily, Kobbi Nissim, Adam D. Smith, Thomas Steinke, Uri Stemmer, and Jonathan Ullman. In STOC, 2016.
- [DFH+15a] Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toni Pitassi, Omer Reingold, and Aaron Roth. In NIPS. 2015.
- [DFH+15b] Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Leon Roth. In STOC, 2015.
- [DMNS06] Cynthia Dwork, Frank Mcsherry, Kobbi Nissim, Adam Smith. In TCC, 2006.
- [HU14] Moritz Hardt and Jonathan Ullman. In FOCS, 2014.
- [RZ16] Daniel Russo and James Zou. In AISTATS, 2016.