

Privacy: the setting of this course

Data providers
ex: individuals



Sensitive data
(HIV status)

Data analyst
ex: health institute,
estimate HIV



Privacy: the objective

Data providers
Privacy protections



Sensitive data
(HIV status)

Data analyst
**Accurate,
informative data
analysis**



Privacy: attempt 1

Learn **nothing** about any specific individual whose data we use in our computation

- I.e., perfect privacy protection; no information *at all* is leaked

Privacy: attempt 1

Learn **nothing** about any specific individual whose data we use in our computation

- I.e., perfect privacy protection; no information *at all* is leaked

But perfect privacy is incompatible with informative data analysis

- Perfect privacy → statistically, no information about each agent contained in computation
- Output of computation independent of data → uninformative and inaccurate

Privacy: attempt 1

Learn **nothing** about any specific individual whose data we use in our computation

- I.e., perfect privacy protection; no information *at all* is leaked

But perfect privacy is incompatible with informative data analysis

- Perfect privacy → statistically, no information about each agent contained in computation
- Output of computation independent of data → uninformative and inaccurate

There is a trade-off between privacy and accuracy

- Need *some* privacy to be leaked for accuracy

Privacy: attempt 2

Learn **almost nothing** about any specific individual whose data we use in our computation

- Now we can look at privacy \leftrightarrow accuracy trade-offs
- Minimize the amount of privacy leaked for a given level of accuracy/usefulness

Privacy: attempt 2

Learn **almost nothing** about any specific individual whose data we use in our computation

- Now we can look at privacy \leftrightarrow accuracy trade-offs
- Minimize the amount of privacy leaked for a given level of accuracy/usefulness

Problem: we can learn from statistical inference

- Imagine I learn from my data analysis that smoking causes cancer
- I am your insurance company and know that you smoke
- ➔ Believe you are at higher risk of concern and can overcharge/harm you!

Privacy: final attempt

Learn almost nothing about any specific individual whose data we use in our computation, **that we would not have learned without their data**

- A single agent's data has *little* effect on what I learn
- What I learn is almost independent of a single agent's data → hard to reconstruct an agent's data from output

Privacy: final attempt

Learn almost nothing about any specific individual whose data we use in our computation, **that we would not have learned without their data**

- A single agent's data has *little* effect on what I learn
- What I learn is almost independent of a single agent's data → hard to reconstruct an agent's data from output

Protects from *additional* harm from giving data, not from *all* harms:

- Statistical inference: smoking → cancer, certain zipcodes → certain ethnicities
- But this is unavoidable:
 - Any such accurate statistic can cause harm
 - Not a property of differential privacy, but rather a privacy impossibility

Differential privacy, informally

- Database D made of agents' sensitive data; each row \Leftrightarrow data of a single agent
- Want to answer query q on database D , without learning any row/agent data

Database D

Name	Has HIV?
Juba	No
Rick	Yes
Homer	No

$q(D)$ = fraction of pop.
 D that has cancer?

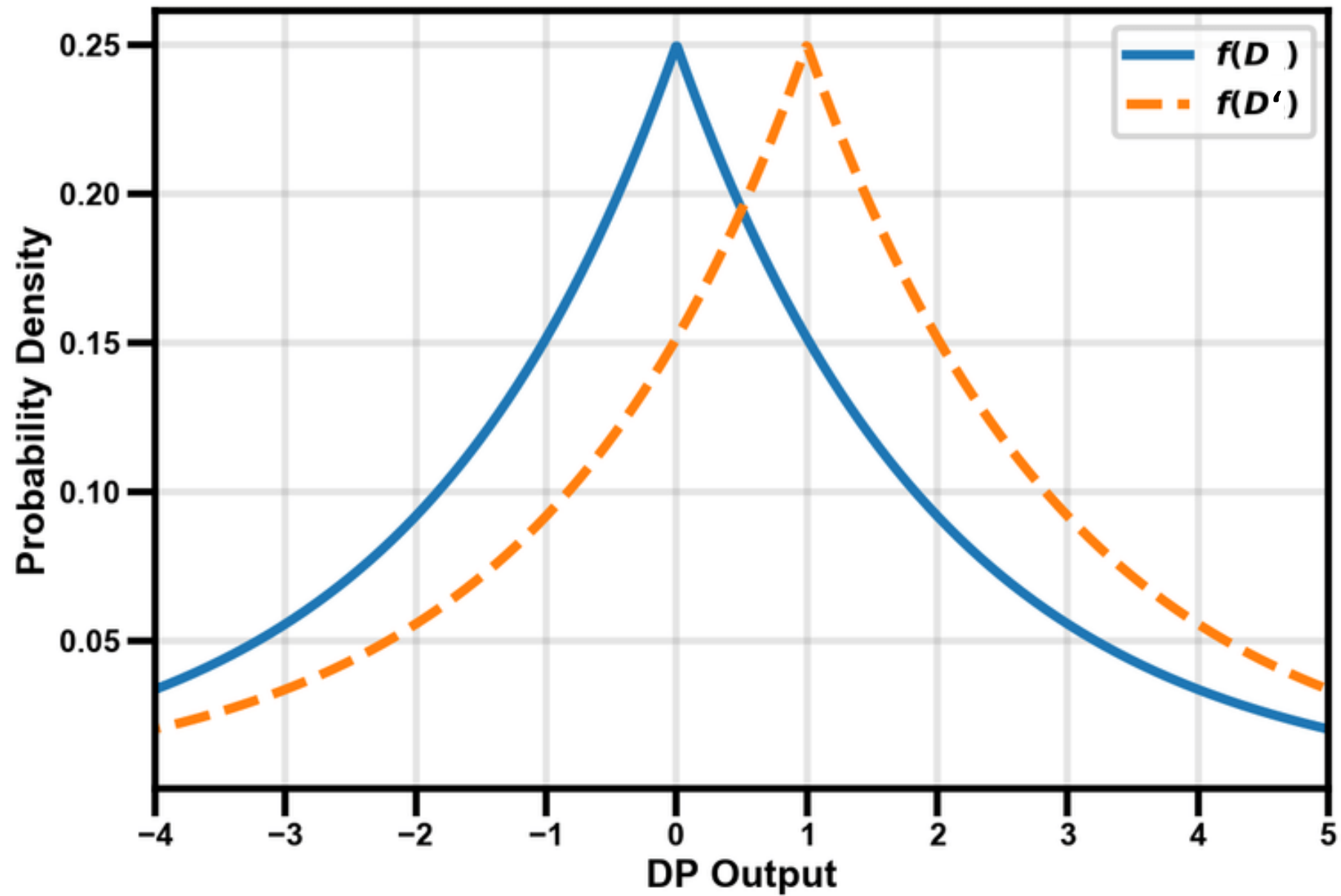
Differential privacy, informally

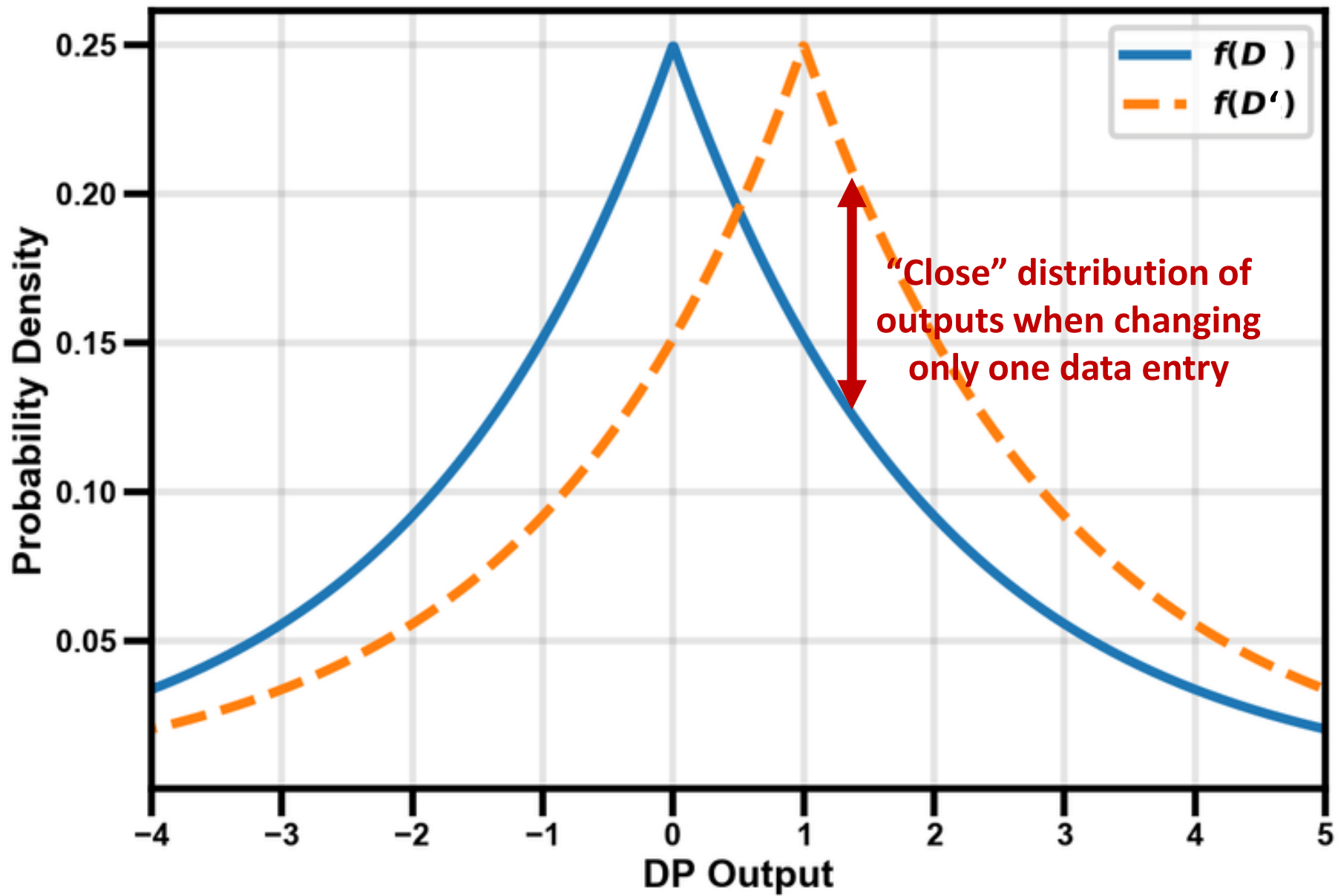
DP asks what happens if change data of single agent?

Database D			Neighboring database D'	
Name	Has HIV?		Name	Has HIV?
Juba	No	→	Juba	No
Rick	Yes		Rick	No
Homer	No		Homer	No

Distribution of outputs of computation *almost unchanged*

- Privacy guarantee: outcome (almost) independent of a single agent's data
- How? Answer $q(D)$ + add **random noise** (from well-chosen dist.)





Differential privacy \neq security

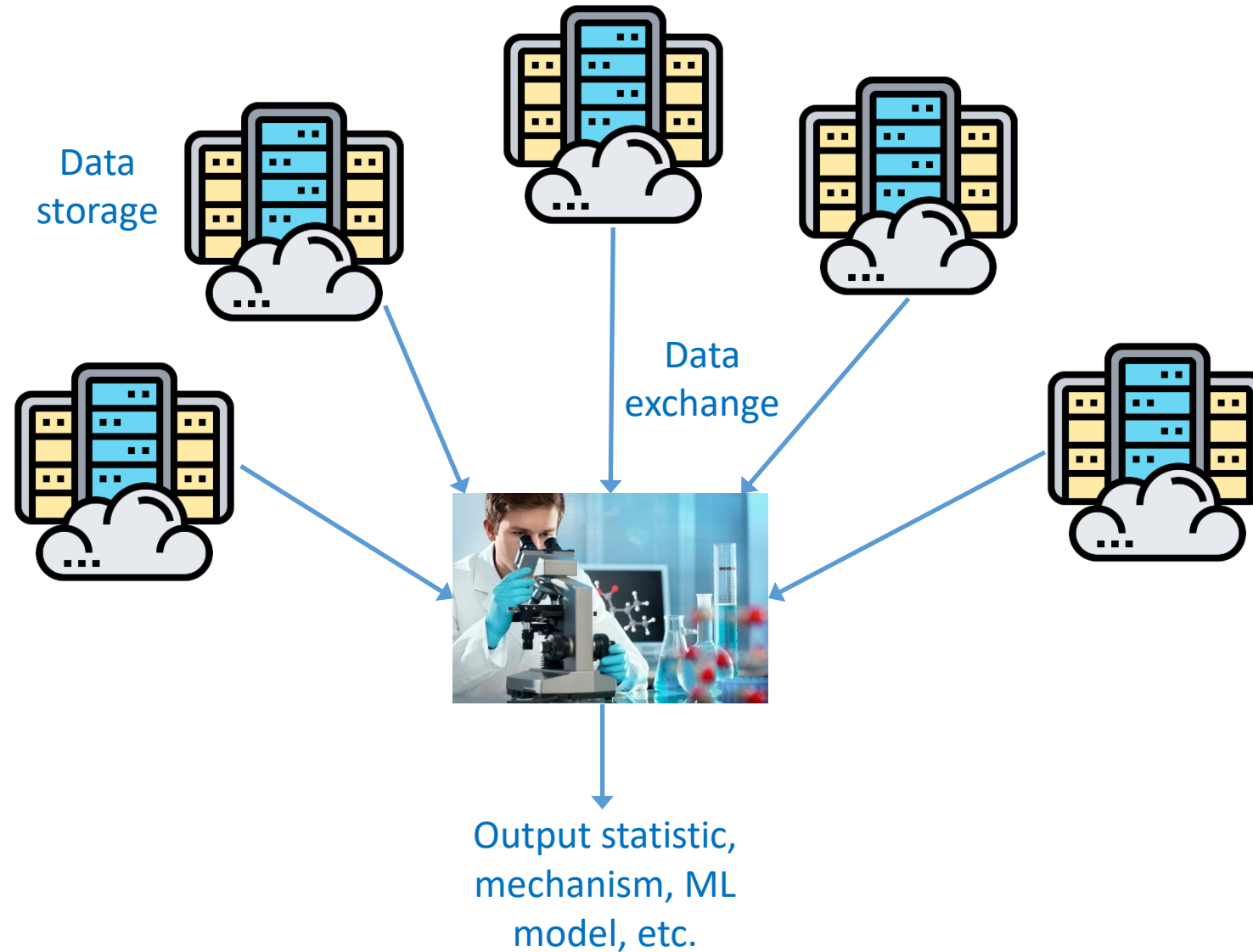
Security and cryptography:

- How can I prevent an attacker from accessing my database?
- How do I protect content of communication between different parties?
- Prevents access to the **raw** data

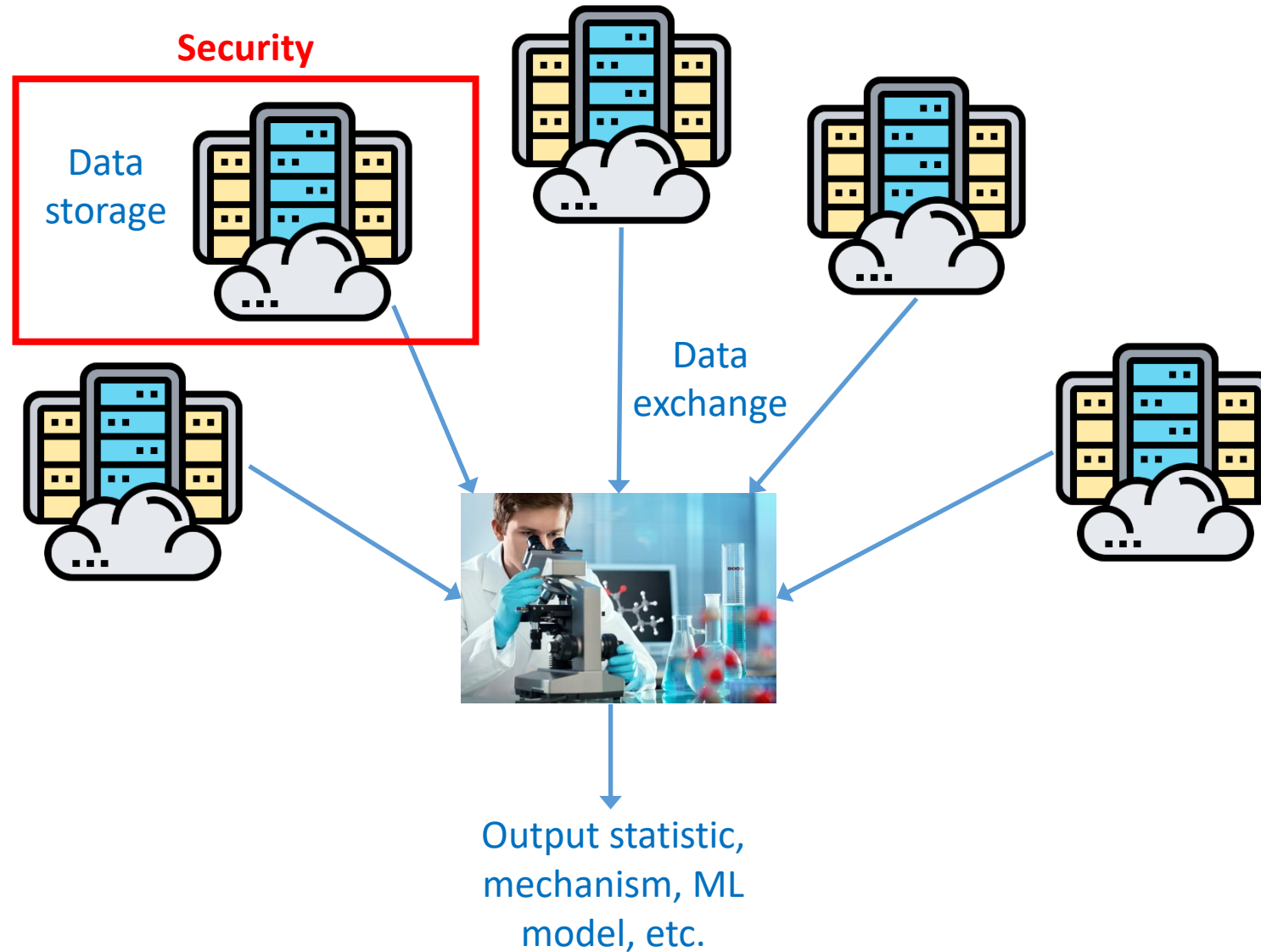
Differential privacy:

- What can I infer from the **output** of my computation?
- Your participation in the algorithm does not change the **output** of the mechanism much \rightarrow cannot infer your data from the output statistic/model/etc.

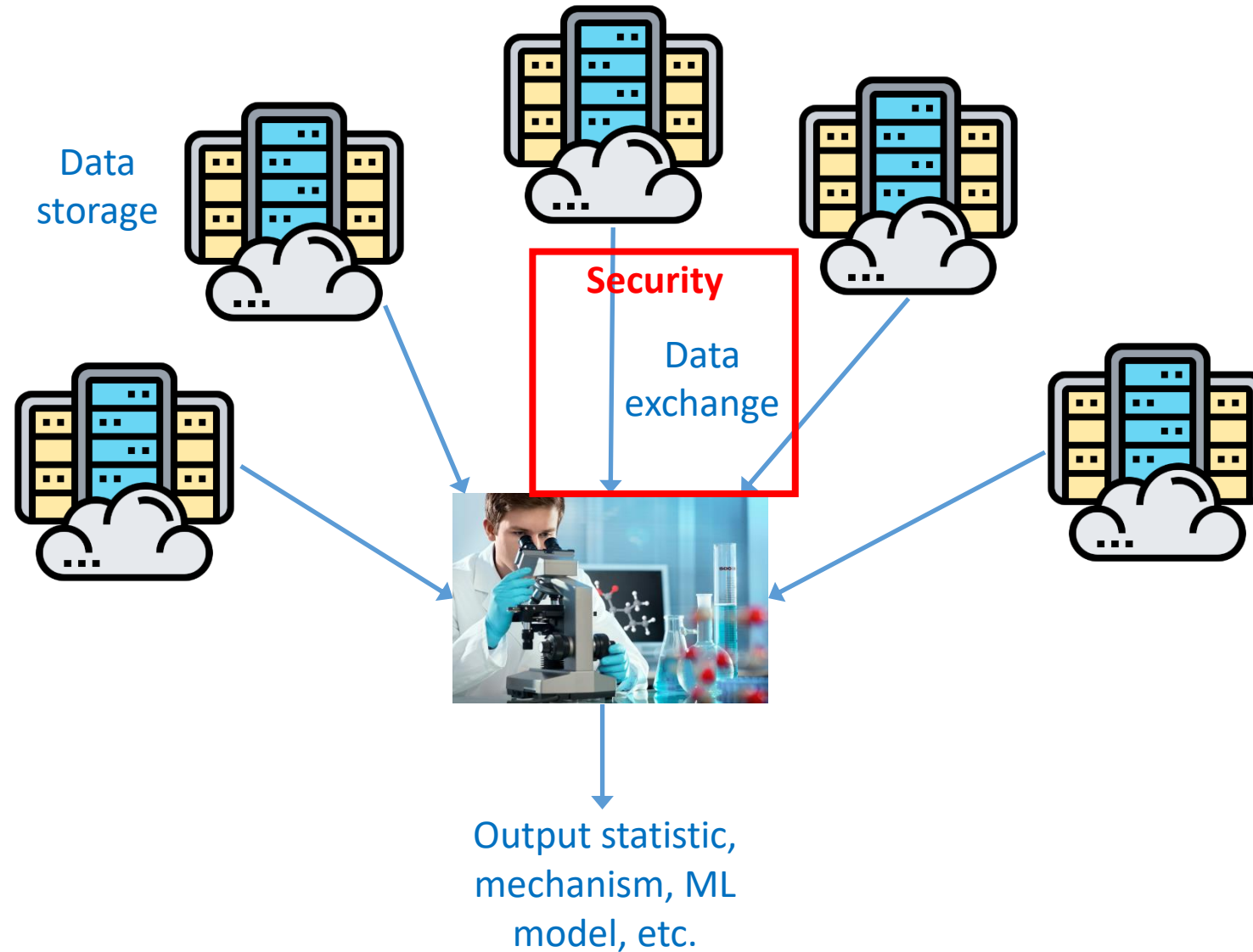
Differential privacy \neq security



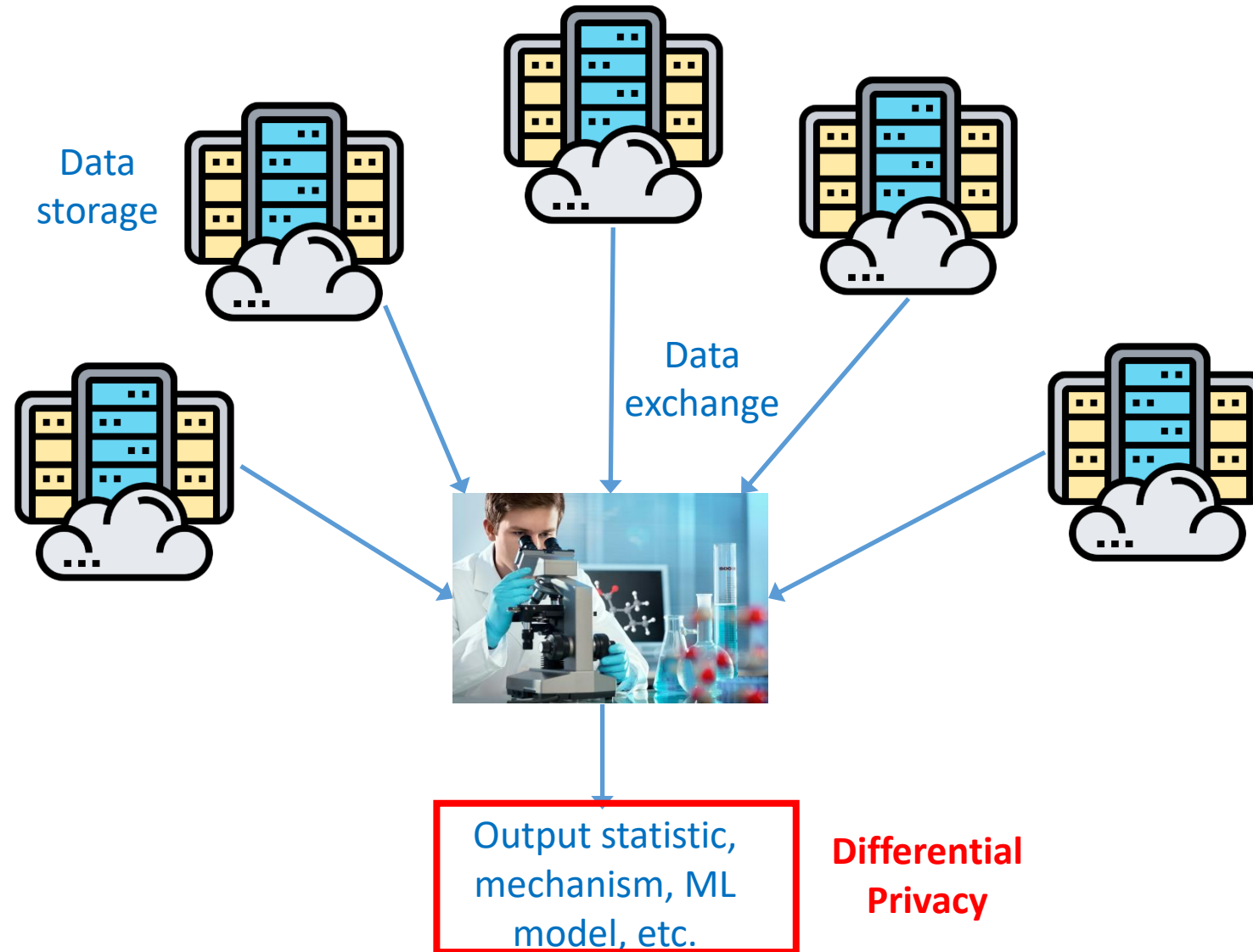
Differential privacy \neq security



Differential privacy \neq security



Differential privacy \neq security



Formal definition of DP