# Controlling privacy risk in database studies for human subject protection(HSP) via a privacy budgeting system

### Hye-Chung Kum[1,2,3], Cason Schmit[1], Alva O. Ferdinand[1,] Theodoros Giannouchos[1]

[1] Texas A&M University Health Science Center, Department of Health Policy & Management, [2] Texas A&M University, Department of Computer Science & Engineering
[3] Population Informatics Research Group (https://research.tamhsc.edu/pinformatics/)

## Background

- Recent advances in information privacy have shown the need to think about ethical data analysis as a budget-constrained problem.
- The goal is to quantify and contain privacy risk under a fixed privacy budget appropriate for the purpose.
- A well-balanced system would allow maximized benefit while maintaining research risks under a fixed budget.
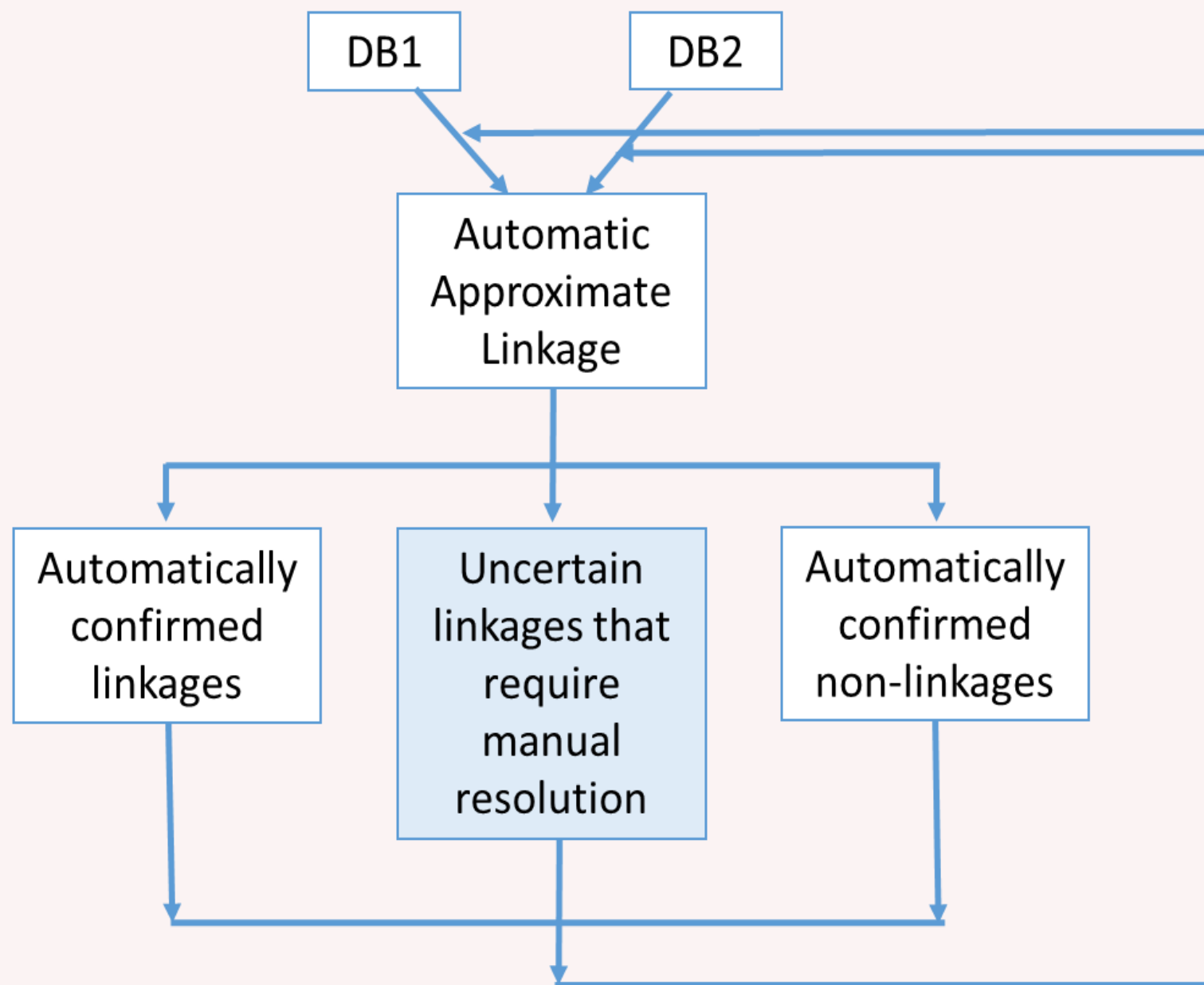


Figure 1 - Approximate RL Human-computer system

## Methods

- We propose a privacy budget system for human subject protection that uses anonymity-set size to define allowable risk in database studies.
- The anonymity-set size, n, is the number of people in the population that share the same identifying information disclosed during research which can quantify the risk of identity disclosure.
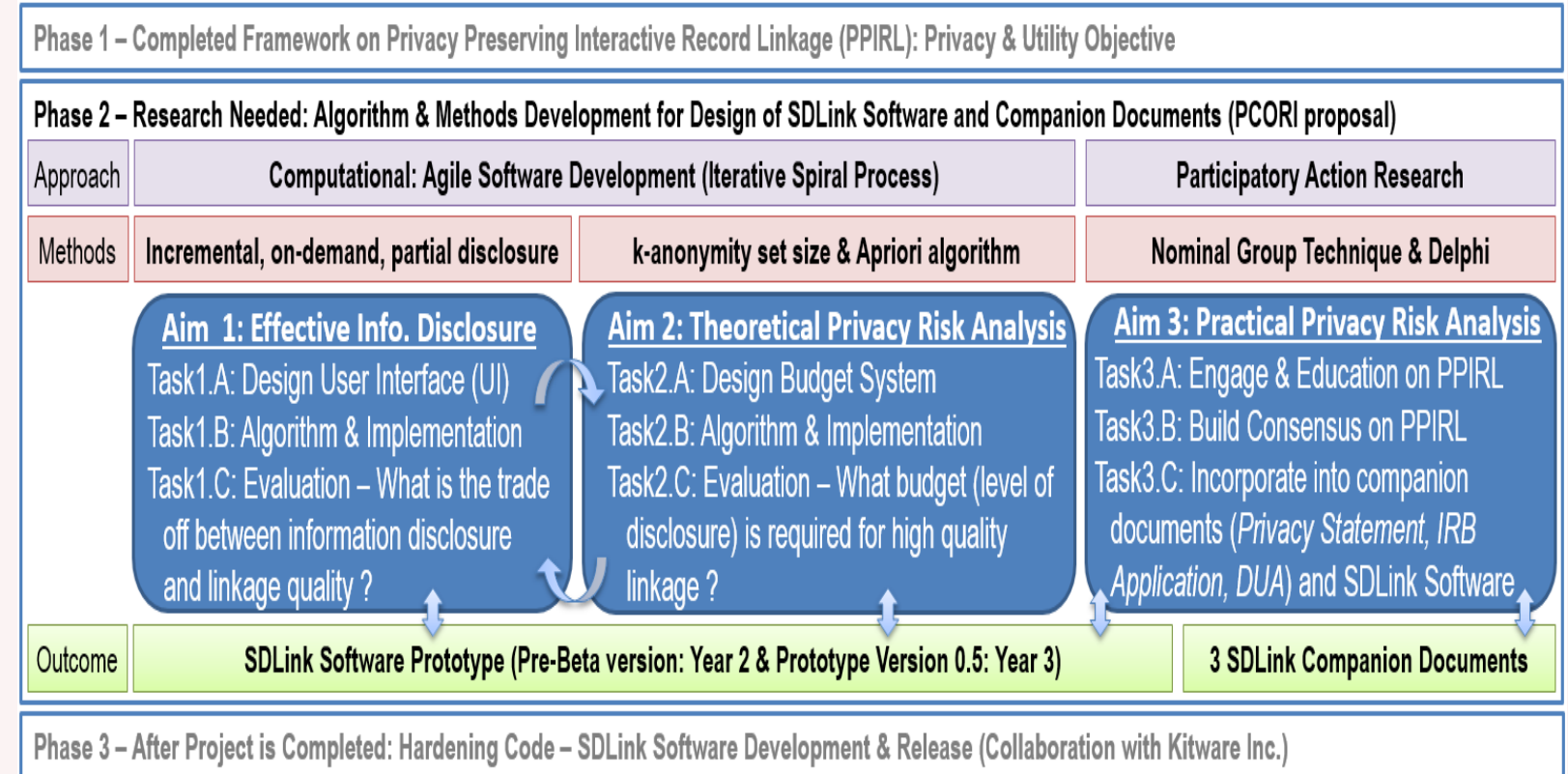


Figure 2 – PPIRL Software Development Phases

## Discussion

- Traditional mechanisms for HSP (e.g., meaningful informed consent, de-identification) improperly prioritize the ethical principle of autonomy by restricting research that advances the principles of beneficence and justice.
- In recognition, recent changes to the Common Rule better facilitate database studies by simplifying secondary data analysis HSP requirements. This privacy budget system aligns with the new Common Rule and promotes a better understanding and quantifying of specific privacy risks in research.
- Novel methods for HSP based on digital technology that can facilitate better reasoning, communication, and negotiations of privacy risk in database studies can improve transparency in database studies ultimately improving human subject protection.



Figure 3 – Disclosure Modes

# Controlling privacy risk in database studies for human subject protection(HSP) via a privacy budgeting system

**Hye-Chung Kum[1,2,3], Cason Schmit[1], Alva O. Ferdinand[1,] Theodoros Giannouchos[1]**

[1] Texas A&M University Health Science Center, Department of Health Policy & Management, [2] Texas A&M University, Department of Computer Science & Engineering
[3] Population Informatics Research Group (https://research.tamhsc.edu/pinformatics/)

## Results



*Scores in each mode*

84.8%   84.1%   84.5%   78.1%   74.5%

100 %   100%   30%   7%   0%

- We can get <u>comparable results to full mode with only 30% disclosure</u> with appropriate masks (moderate mode)
- As we mask more values for privacy, <u>quality of results start to suffer</u> (p<0.001)
- However, even legally de-identified data with proper masks can be linked properly for most situations
  - 0% disclosure still had 75% accuracy
  - Incremental disclosure can significantly improve privacy protection with negligible impact on quality of linkage

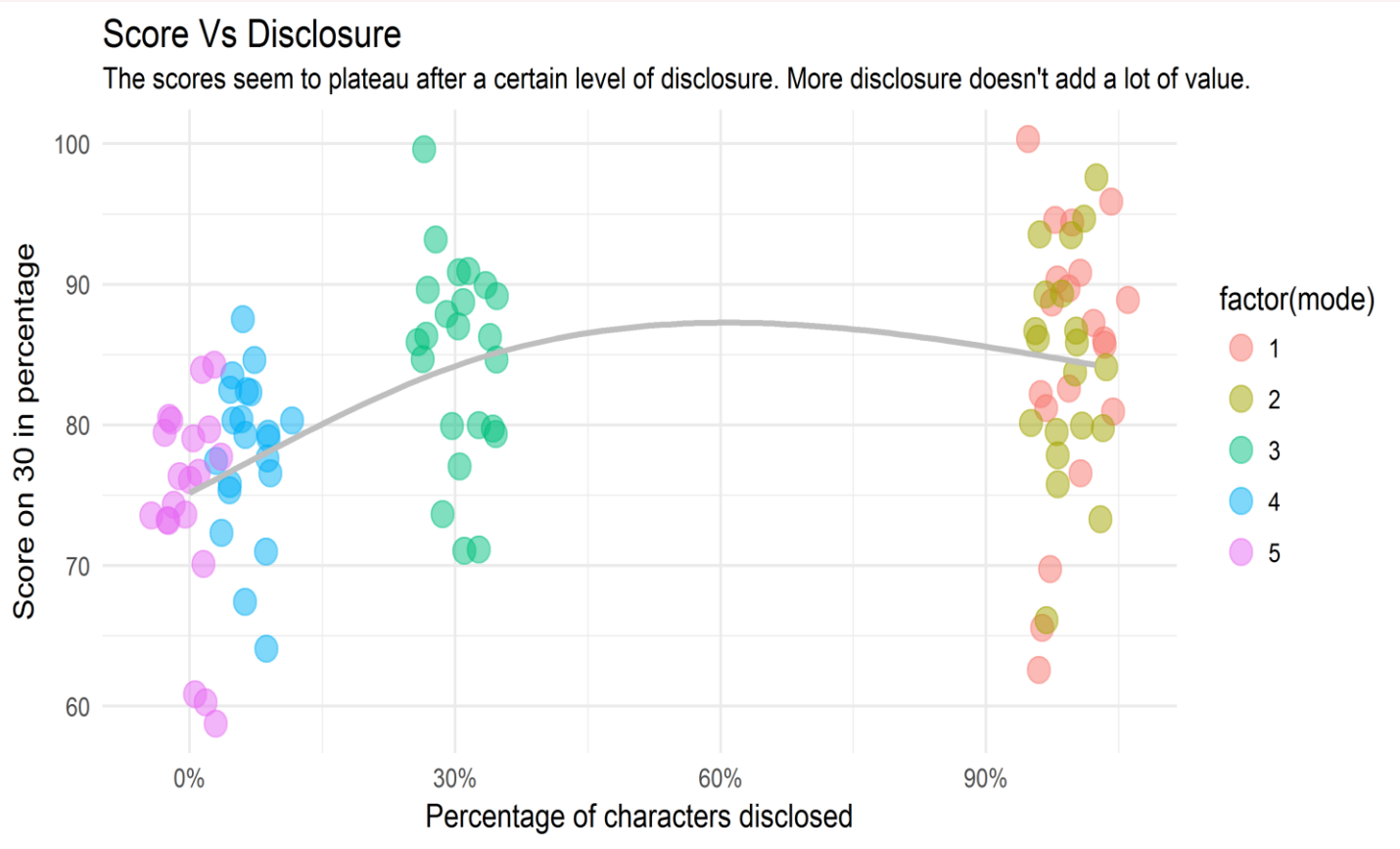### Figure 4 – Accuracy Score by Disclosure Mode



Score Vs Disclosure
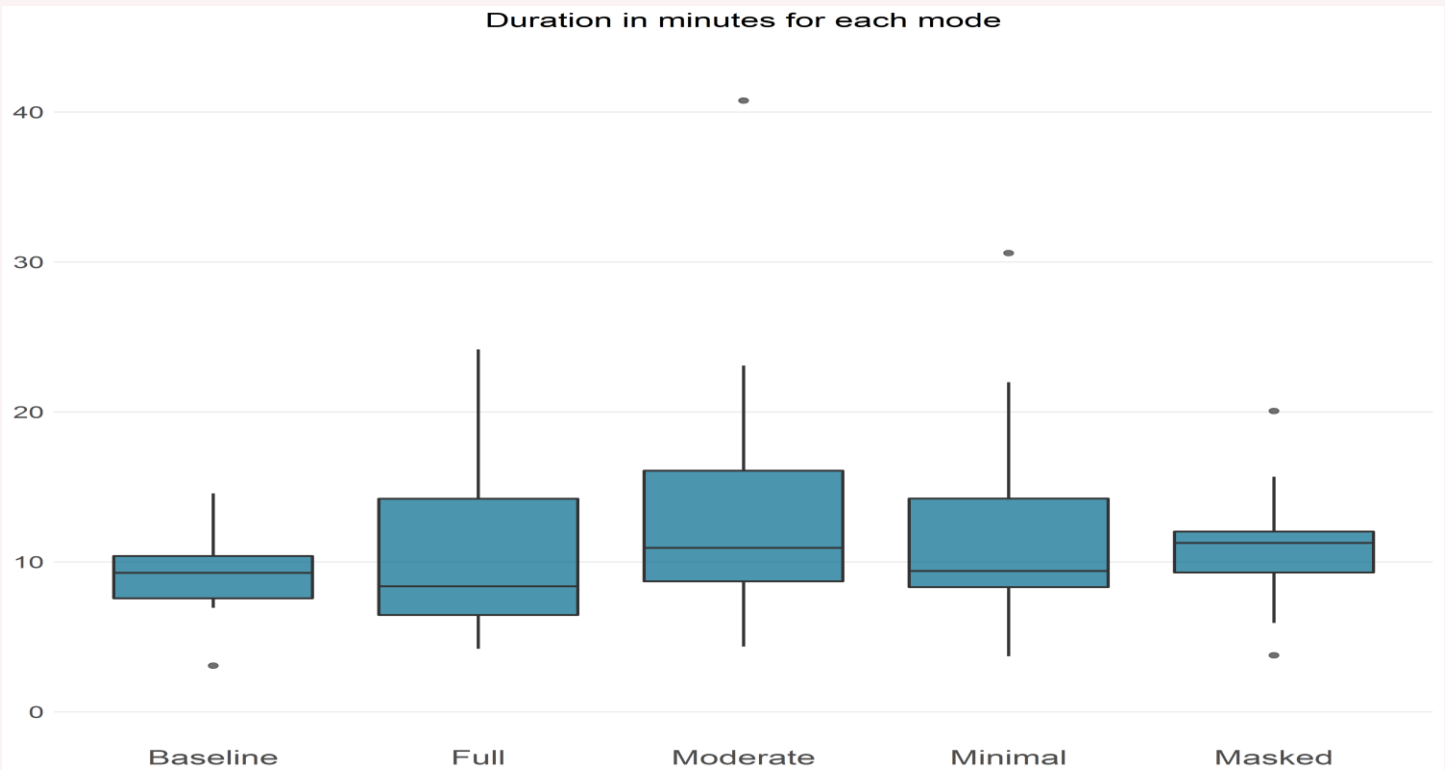The scores seem to plateau after a certain level of disclosure. More disclosure doesn't add a lot of value.

### Figure 5 –Score vs Distance



*Duration in minutes for each mode*

## Conclusion

- Mathematically, the <u>identity disclosure risk</u> is <u>inversely related to the number of entities</u> in the population that share the information disclosed during research.
- If the <u>information refers to one and only one person</u> in the population, then the <u>identity</u> of the person has been <u>fully disclosed</u> by the information revealed.
- On the other hand, <u>if the information disclosed is identical for multiple people</u> (=n), referred to as the <u>anonymity-set size</u>, then the information is <u>less revealing</u> as it could refer to any one of the n people.
- For example, when <u>a frequently occurring name</u> (e.g., "Mary") is disclosed, there is <u>low probability that the identity</u> of a specific person named Mary is <u>revealed</u>. In comparison, when a <u>rare name</u> (e.g., "Hye-Chung") is disclosed, there is a <u>greater risk that the identity will be ascertained.</u>
- <u>By combining all anonymity-set sizes</u> for different information disclosed during a given study, <u>we design a privacy risk score based on the background population</u>.
- When the background population information is not readily available, as is often the case, the study population can be used to calculate the most conservative privacy risk score.
- The proposed privacy budget system can enable all stakeholders (e.g., database researchers, IRBs, the public) to (1) <u>accurately reason about</u> (2) <u>communicate on</u>, and (3) <u>agree on the acceptable quantified privacy risk</u> considering the anticipated benefits of a given study.
- Ideally, the budget would be a simple measurement on the probability of identity disclosure for the population. It will also allow <u>for better monitoring of risk to increase transparency</u>.

## Limitations

- The proposed method has <u>not been implemented</u> and no evaluation is currently planned.
- More time is required to finalize and evaluate the concept.
- This poster provides a great opportunity to gather wide input from the IRB community on the novel concept before the privacy budget design is finalized and implemented.

### Figure 6 –Time by Disclosure Method