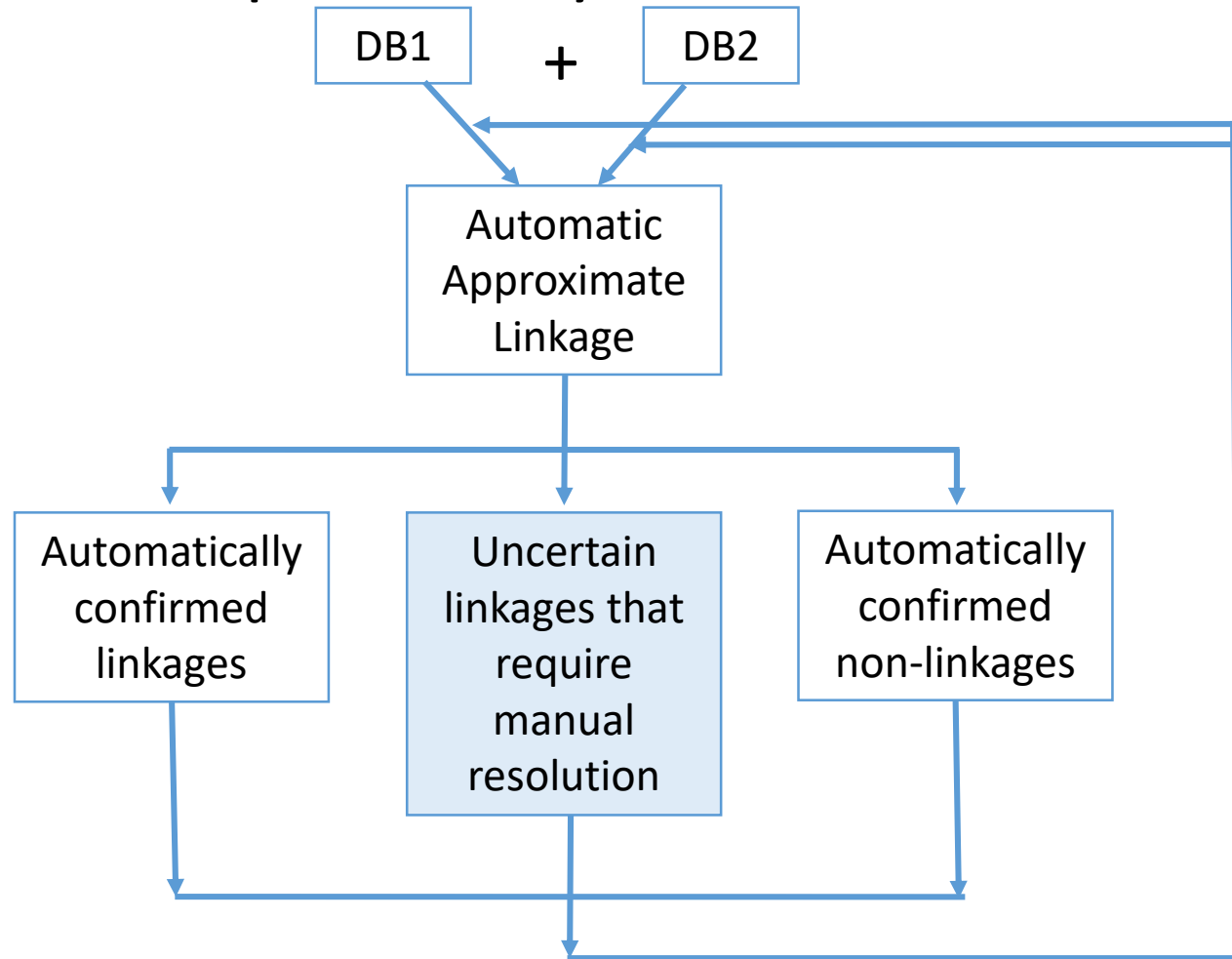


# Approximate RL Human-computer system



- Standardize
- Clean Data
- Build Training Data

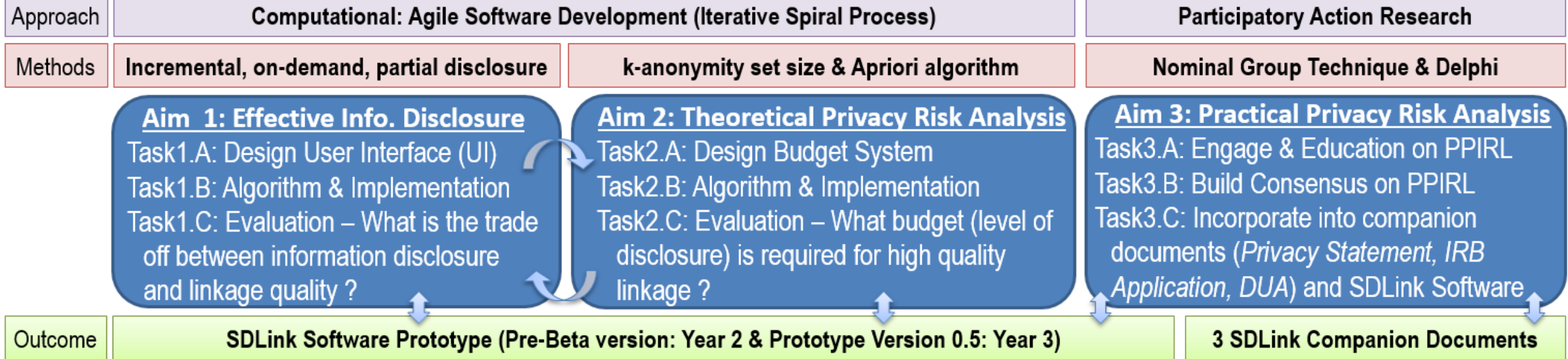
- 75%-80% automatics
- 15%-25% manual resolution

# Privacy Preserving Interactive Record Linkage (PPIRL) via Information Suppression

<https://pinformatics.org/ppirl/>

**Phase 1 – Completed Framework on Privacy Preserving Interactive Record Linkage (PPIRL): Privacy & Utility Objective**

**Phase 2 – Research Needed: Algorithm & Methods Development for Design of SDLink Software and Companion Documents (PCORI proposal)**



**Phase 3 – After Project is Completed: Hardening Code – SDLink Software Development & Release (Collaboration with Kitware Inc.)**

# Aim 1: Effective Information Disclosure

- July 2017: User Study 1
  - The study had a total of 104 participants
    - ~20 participants for each of the five modes
    - There were 61 males and 42 females, and one participant did not specify gender.
    - Ages ranged from 18 to 43 years, and the median age was 24 years.
    - About 65% of the participants were from the United States and had English as their native language.
    - About 57% of the participants were either pursuing or already had a graduate degree, and the remaining participants were undergraduate university students.
  - 30 questions

# BASE mode

Pair	ID	First name	Last name	DoB (M/D/Y)	Sex	Race
1	8000002767	JUDE	WILLIAM	09/09/1906	M	W
	8000003567	JUDE	WILLIAM JR	09/09/1960	M	B

# FULL mode (Icons & colors)

Pair	ID	FFreq	First name	Last name	LFreq	DoB (M/D/Y)	Sex	Race
1	8000002767	①	JUDE	WILLIAM	①	09/09/1906	M	W
	8000003567	①	JUDE	WILLIAM JR	①	09/09/1960	M	B

# MODERATE mode (Close same & partial IDs)

Pair	ID	FFreq	First name	Last name	LFreq	DoB (M/D/Y)	Sex	Race
1	*****27**	①	✓	WILLIAM	①	09/09/1906	M	W
	*****35**	①	✓	WILLIAM JR	①	09/09/1960	M	B

# MINIMUM mode (Partial names, dates etc)

Pair	ID	FFreq	First name	Last name	LFreq	DoB (M/D/Y)	Sex	Race
1	*****27**	①	✓	*****	①	**/**/**06	M	@
	*****35**	①	✓	***** JR	①	**/**/**60	M	&

# MASKED mode (Only symbols)

Pair	ID	FFreq	First name	Last name	LFreq	DoB (M/D/Y)	Sex	Race
1	*****@**	①	✓	*****	①	**/**/**@	✓	@
	*****&**	①	✓	***** &	①	**/**/**&	✓	&

# ENCRYPTED mode (Encrypt Data)

Pair	ID	FFreq	First name	Last name	LFreq	DoB (M/D/Y)	Sex	Race
1	@zz@@@e@@@yc@tcflg==	①	HOMwdz8KpFKaTfPE+qr8Xw==	✓	①	/KSKzJ2U5C/fpHmkMqZPqw==	~	%
	&qw&p&&&&m&&&v1m&==	①	o1fSci26GzxKx41n11kRuQ==	✓	①	bJupC1Skjj/bmw9DRq07vw==	~	^

# Mode 2: Full Mode

## Full disclosure with markup

Pair	ID	FFreq	First name	Last name	LFreq	DoB(M/D/Y)	Sex	Race
1	8000002767 ✘	①	JUDE	WILLIAM	①	09/09/1906	M	W
	8000003567	①	JUDE	WILLIAM JR +	①	09/09/1960	M	DIFF B
2	0000006947	①	BRYANT	MADELINE	①	05/02/1962	F	W
	0000006947	25	MADELINE	BRYANT	...	05/02/1962	F	W
3	9000018540 DIFF	...	SALLY	BYRD	...	07/04/1960	F	W
	6000008928	∞	JOHN	BYRD	...	04/07/1960	DIFF M	?

# Mode 3: Moderate Mode

## Moderate disclosure with markup

Pair	ID	FFreq	First name	Last name	LFreq	DoB(M/D/Y)	Sex	Race
1	*****27** ✘	①	✓	WILLIAM	①	09/09/1906	M	W DIFF
	*****35**	①	✓	WILLIAM JR	①	09/09/1960	M	B
2	✓	①	&&&&&	↔ ↔	①	✓	F	✓
	✓	2-5	@@@@@@@@	↔ ↔	...	✓	F	✓
3	@@@@@@@@@@ DIFF	...	SALLY	✓	...	07/04/1960	F	*
	&&&&&&&&&&	∞	JOHN	✓	...	04/07/1960	M	?

# Mode 4: Minimum Mode

## Minimum disclosure with markup

Pair	ID	FFreq	First name	Last name	LFreq	DoB (M/D/Y)	Sex	Race
1	*****27** ✘	①	✓	*****	①	**/**/**06	M	@
	*****35**	①	✓	***** + JR	①	**/**/**60	M	(DIFF) &
2	✓	①	&&&&&	↔ ↔	@@@@@@@@	✓	F	✓
	✓	2-5	@@@@@@@@	↔ ↔	&&&&&	...	✓	F
3	@@@@@@@@@@ (DIFF)	...	@@@@@	✓	...	07/04/****	F	*
	&&&&&&&&&&&&	∞	(DIFF) &&&&	✓	...	✘ 04/07/****	(DIFF) M	? 

# Mode 5: Masked Mode


## Masked disclosure with markup

Pair	ID	FFreq	First name	Last name	LFreq	DoB (M/D/Y)	Sex	Race
1	*****@** X	①	✓	*****	①	**/**/**@	✓	@
	*****&*	①	✓	***** +	①	**/**/**& ⇌	✓	& DIFF
2	✓	①	&&&&&	↔	@@@@@@@@	①	✓	✓
	✓	2-5	@@@@@@@@	↔	&&&&&	...	✓	✓
3	@@@@@@@@@@ DIFF	...	@@@@@	✓	...	@/&/**	@	@
	&&&&&&&&&&&&	∞	&&&	✓	...	&/@/** X	DIFF	?








# Mode 6: Encrypted Mode

## Encrypted disclosure with markup

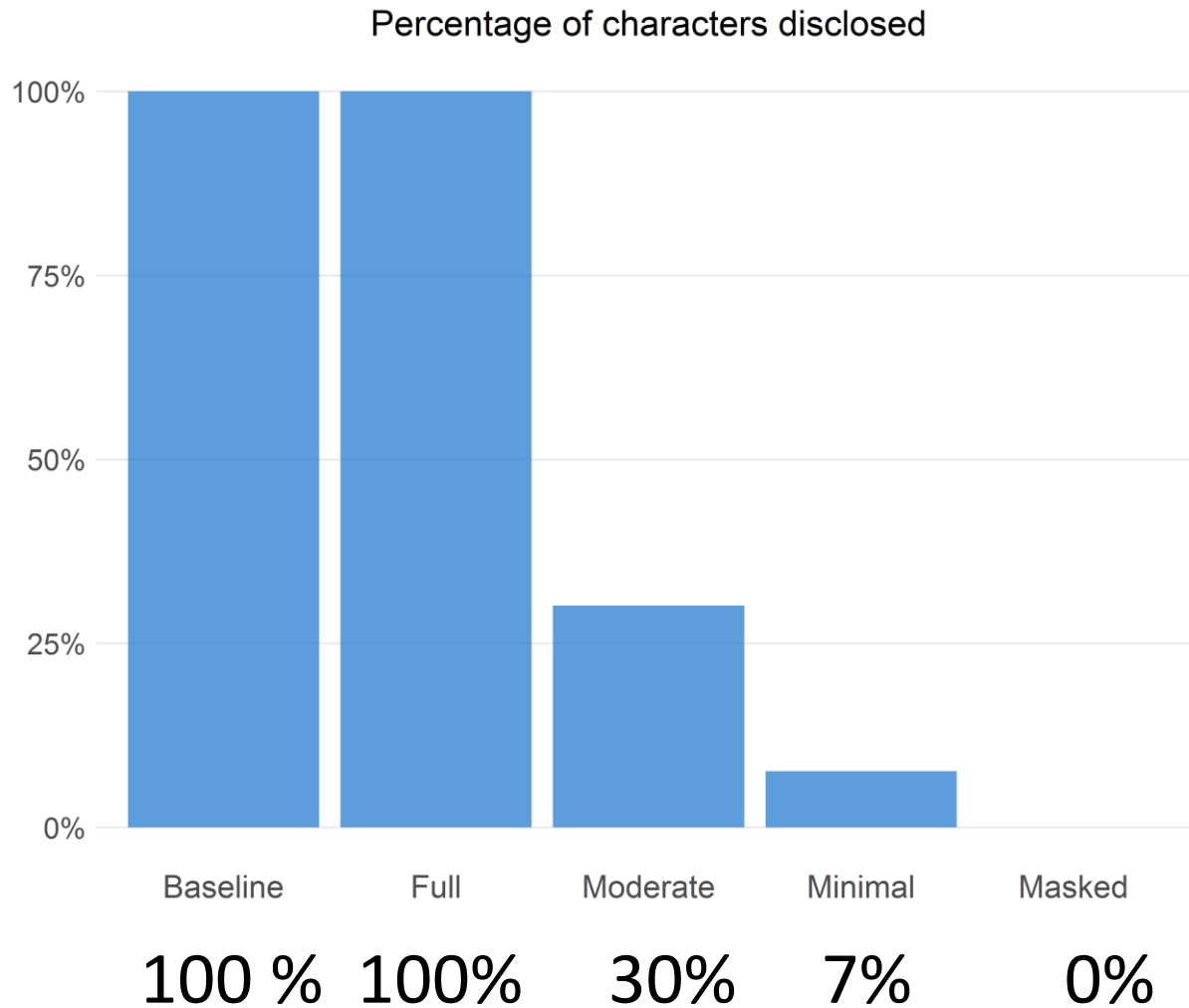
Pair	ID	FFreq	First name	Last name	LFreq	DoB (M/D/Y)	Sex	Race	
31	@@zz@@@e@@@@@yc@tcflg== <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	∞	✓		Pgi+8vEbeh4nP757N9zGdg== <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	①	/KSKzJ2U5C/fpHmkMqZPqw== <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	~	% <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>
	&&qv&p&&&&&m&&&&v1m&&==	①	✓		9SGxq1uytTwBKSC8SpQx8A==	①	bJupC1skjj/bmw9DRq07vw==	~	Λ
32	✓	2-5	#####		#####	...	✓	>	✓
	✓	①	#####	#####	①	✓	>	✓	
33	@@l@i@g@@@os@@bn@@@@g== <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	①	H0Mwdz8KpFKaTfPE+qr8Xw== <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	✓		①	JPHm/tfJf/Sa38z+PthPYQ== <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	>	% <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>
	vx&+&&h&v&+&xnmyqmaa&&==	①	o1fSci26GzxKx41n1lkRuQ==	✓		①	AgsX5d/vZ1tRukT6GTxCZw==	~	?

# Encrypted Mode vs Full Mode

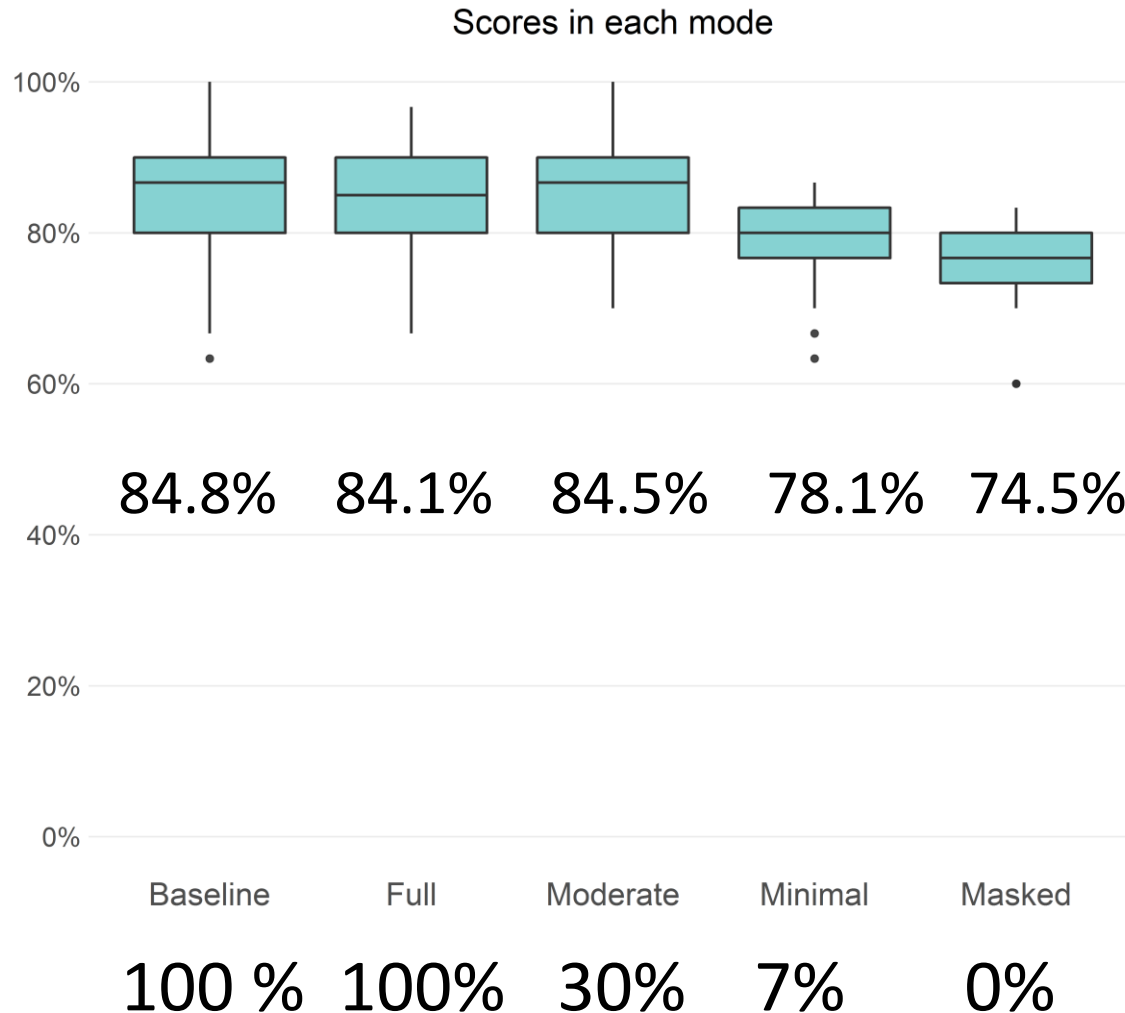
Pair	ID	FFreq	First name	Last name	LFreq	DoB(M/D/Y)	Sex	Race	
31	@@zz@@@e@@@@@yc@tcflg== <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	∞	✓		Pgi+8vEbeh4nP757N9zGdg== <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	①	/KSKzJ2U5C/fpHmkMqZPqw== <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	~	% <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>
	&&qw&p&&&&&m&&&&v7m&&==	①	✓		9SGxq1uytTwBKSC8SpQx8A==	①	bJupC7Skjj/bmw9DRq07vw==	~	^
32	✓	2-5	#####		#####	...	✓	>	✓
	✓	①	#####		#####	①	✓	>	✓
33	@@l@i@q@@@os@@bn@@@@g== <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	①	H0Mwdz8KpFKaTfPE+qr8Xw== <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	✓		①	JPHm/tFJf/Sa38z+PthPYQ== <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	>	% <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>
	vx&+&&h&v&+&xnmymqaa&&==	①	o1fSci26GzxKx41n1kRuQ==	✓		①	AgsX5d/vZ1tRukT6GTxCZw==	~	? <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>

Pair	ID	FFreq	First name	Last name	LFreq	DoB(M/D/Y)	Sex	Race
1	8000002767 ✗	①	JUDE	WILLIAM +	①	09/09/1906 ⇔	M	W <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>
	8000003567	①	JUDE	WILLIAM JR	①	09/09/1960	M	B
2	0000006947	①	BRYANT	 MADELINE	①	05/02/1962	F	W
	0000006947	2-5	MADELINE	 BRYANT	...	05/02/1962	F	W
3	9000018540 <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	...	SALLY <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	BYRD	...	07/04/1960	F <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	W
	6000008928	∞	JOHN	BYRD	...	04/07/1960 	M <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>	? <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">DIFF</span>

# Percentage of Characters Disclosed



# Accuracy Score by Disclosure Mode

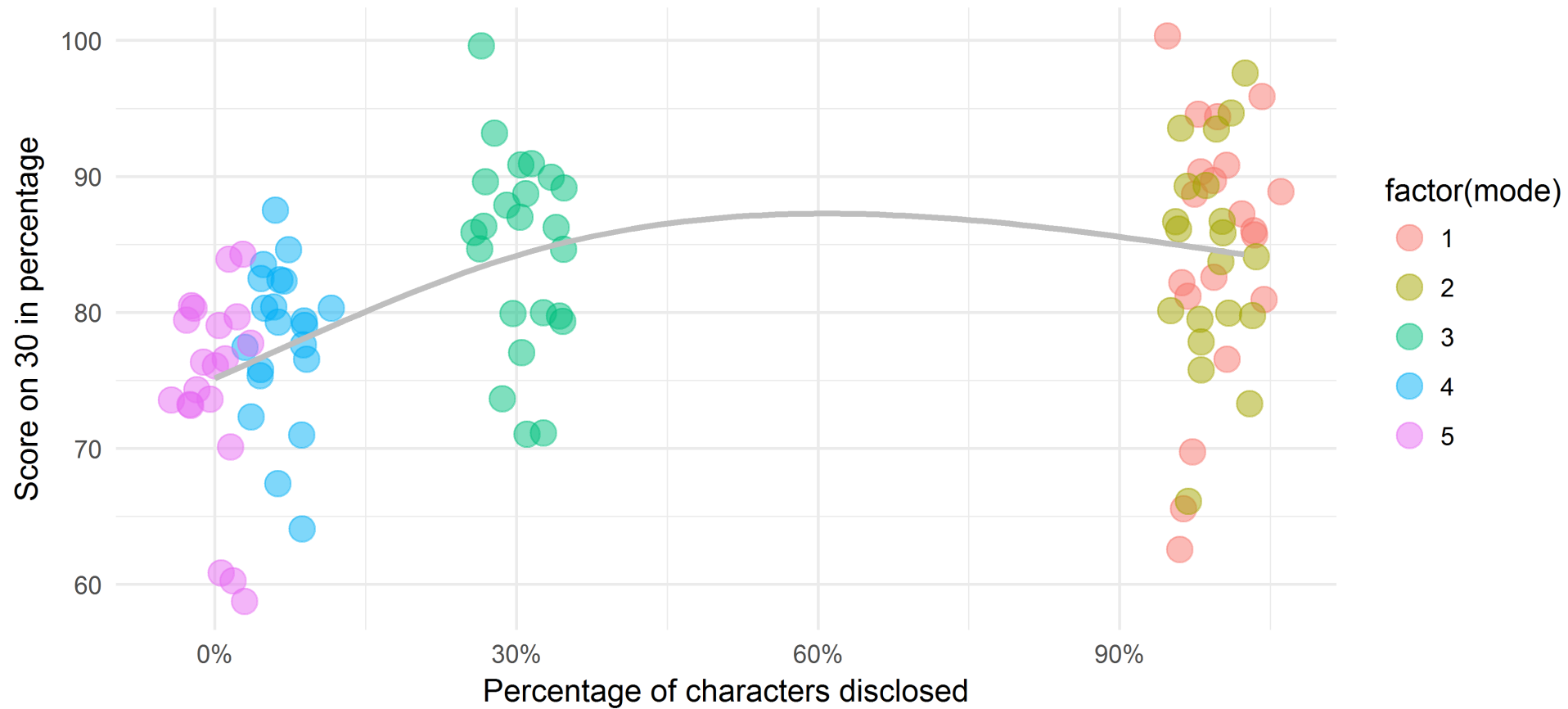


- We can get comparable results to full mode with only 30% disclosure with appropriate masks (moderate mode)
- As we mask more values for privacy, quality of results start to suffer ( $p < 0.001$ )
- However, even legally de-identified data with proper masks can be linked properly for most situations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

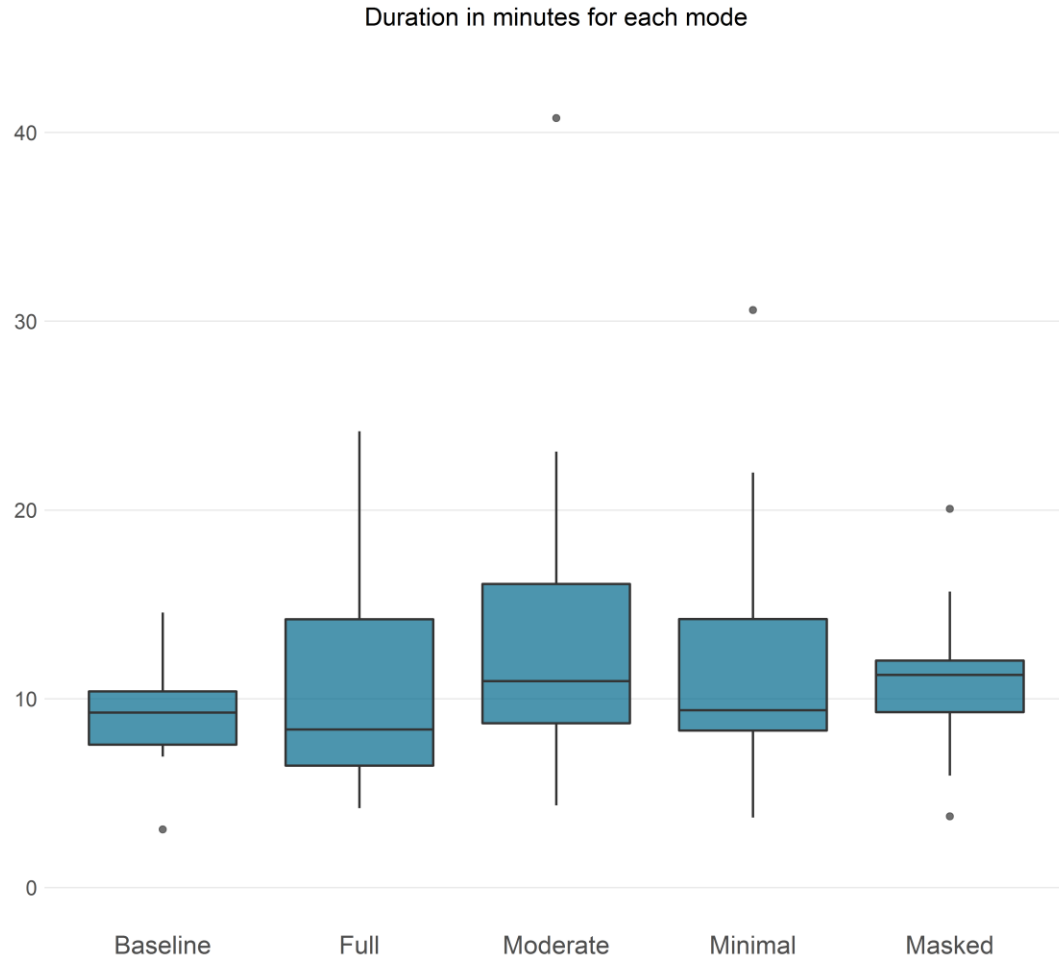
# Score vs Disclosure

## Score Vs Disclosure

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

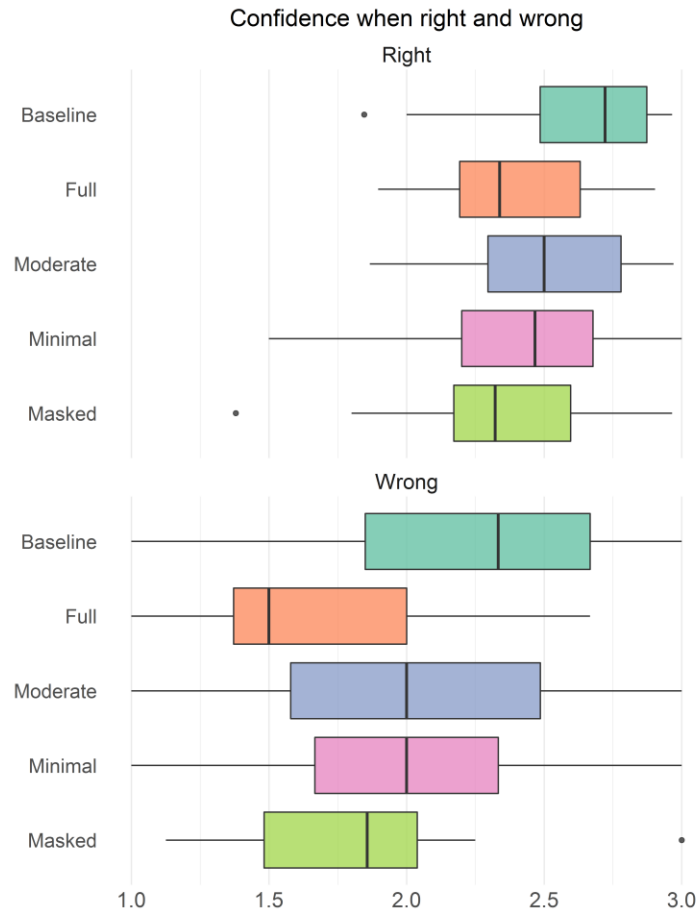


# Time by Disclosure Mode



- Comparable across all modes
- More information (supplemental mark up and frequency icons) has more variability among participants
  - Probably due to differences in participants speed of processing information

# Confidence Level by Correctness of Decision



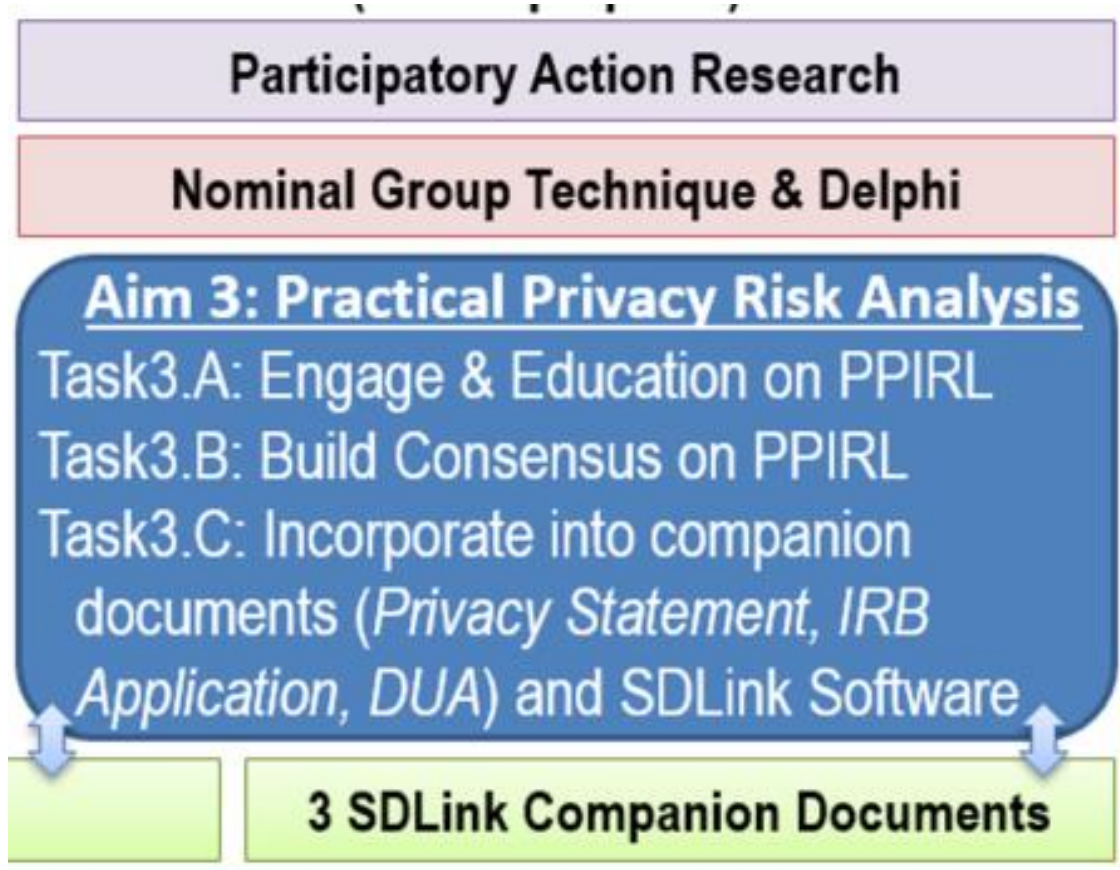
- Higher confidence when answers are correct (top) compared to when answers were wrong (bottom)
- Full mode is least confident when wrong answer
  - More information introduces more uncertainty in wrong decisions, but not sufficient to change the answer

# PPIRL

- Aim 2: Theoretical Privacy Risk Analysis
  - User Study 2: Spring 2018
- Fall 2018: Beta release
  - Summative evaluation: UAB & UT Houston



# Aim 3: Practical Privacy Risk Analysis

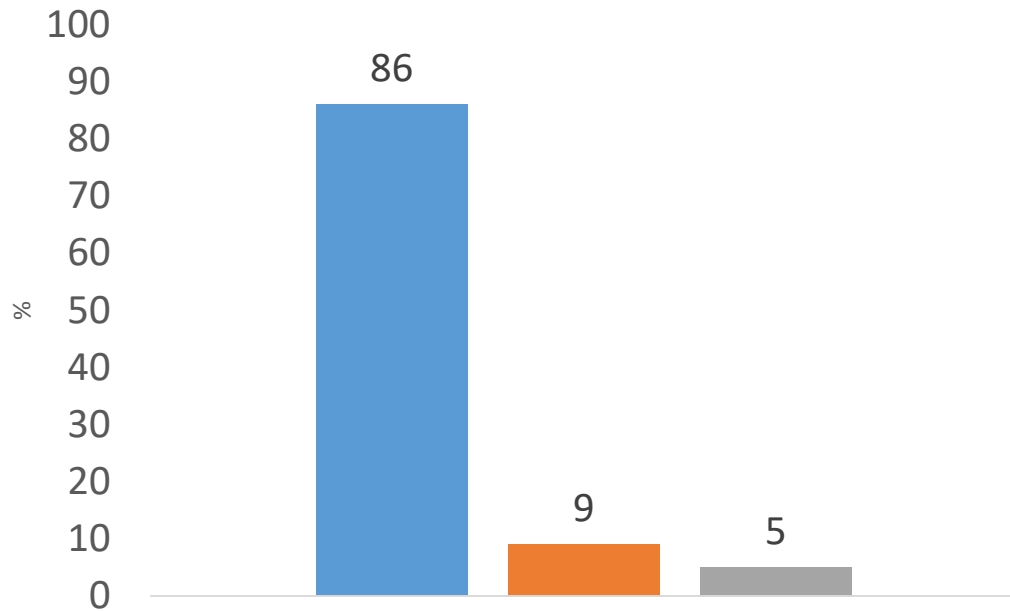


- Template IRB application & DUA
  - Nov 2017: ELSI NGT Session
  - Nov 2018: ELSI Delphi
- Privacy statement
  - Feb 2018: Patient NGT Session
  - Feb 2019: Patient Delphi
  - Feb 2020: Summative Evaluation (Survey)
- User Committee Meeting
  - Every six months: Feb & Aug

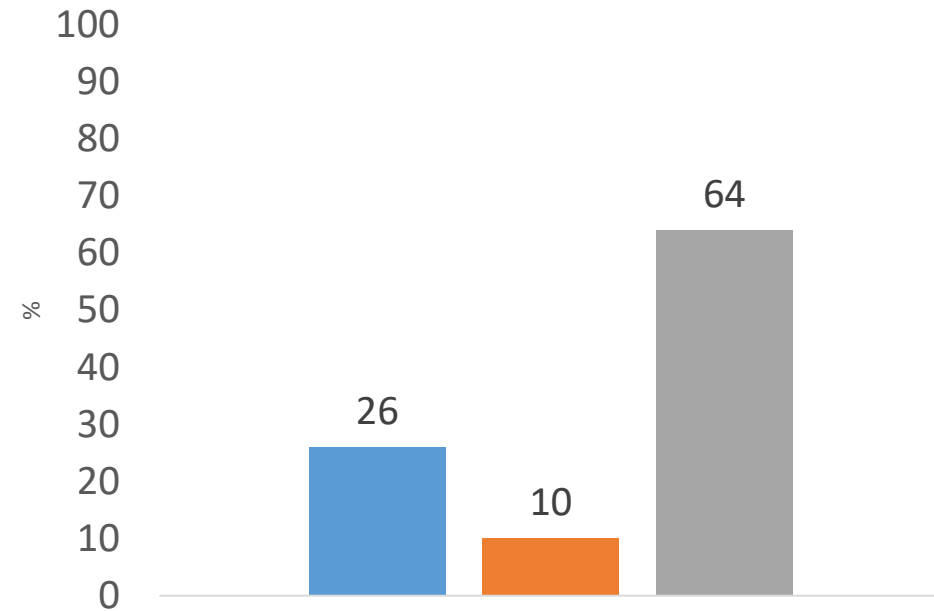
# Improving Methods for Linking Secondary Data Sources for CER/PCOR

- A PCORI Record Linkage Project at Duke University
  - Sean O'Brien & Emily O'Brien
  - July 2014-
- O'Brien E.C., Rodriguez A.M., Kum H.-C., Schanberg L., O'Brien S.M., Setoguchi S. **Patient perspectives on the linkage of health data for clinical research: insights from a survey in the United States.** Oral presentation (#O17-3) at the 2017 World Congress of Epidemiology; Saitama, Japan. August 20, 2017.

# Figure 1. Data Sharing Comfort (n=3516): Sharing PII confidentially



Comfortable with my health data being confidentially shared with researchers, as long as personal information like name and social security number is not available to researchers

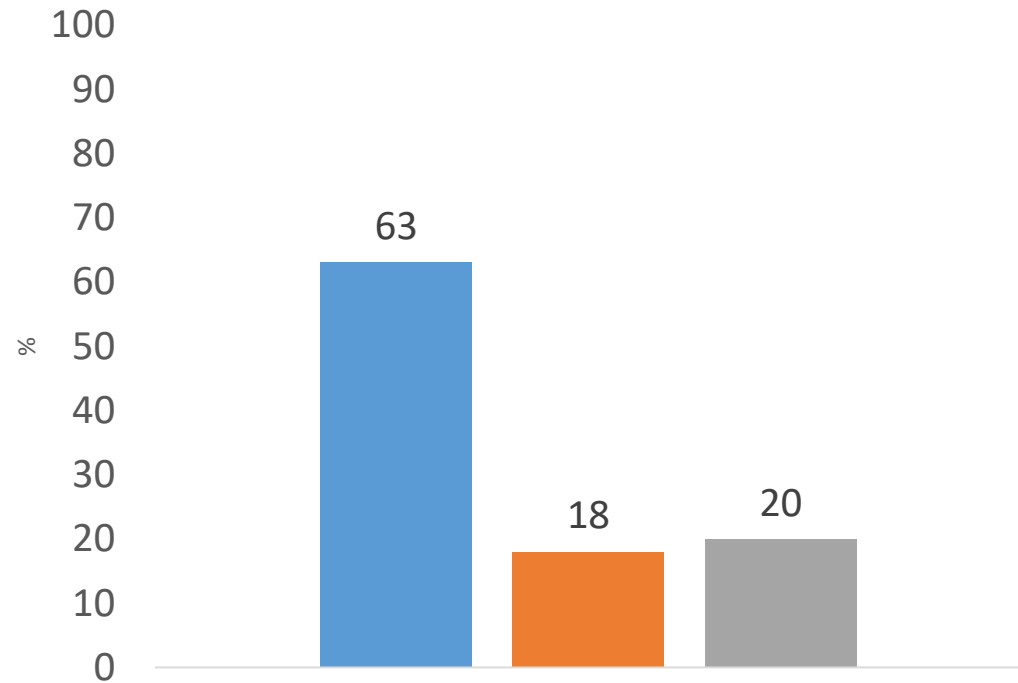


Comfortable with my electronic health data being confidentially shared with health care researchers, EVEN IF personal information like my name and social security number is available

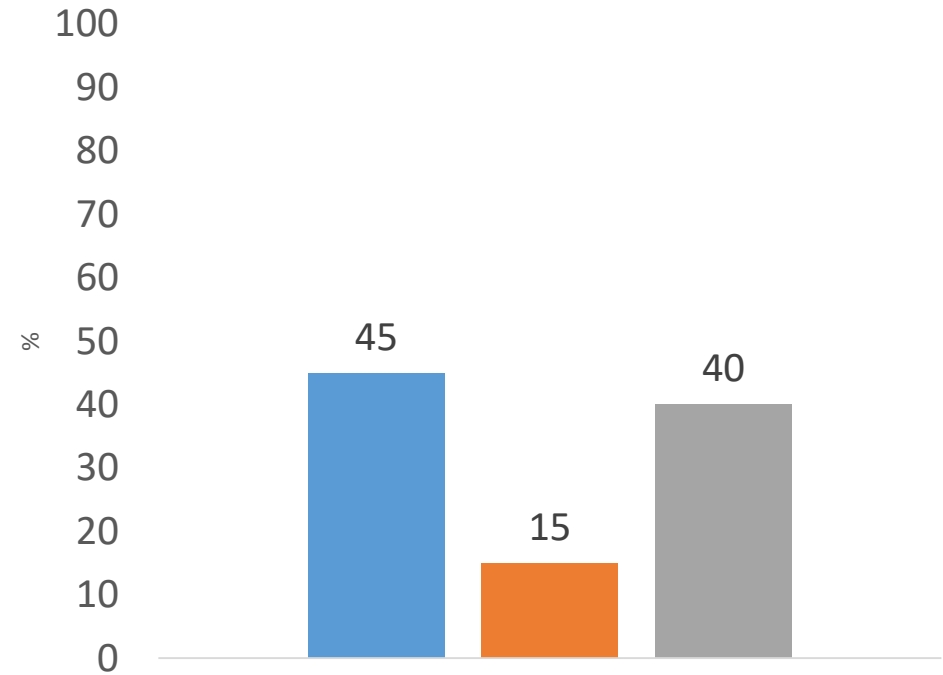
■ Completely/Somewhat Agree ■ Moderately Agree ■ Completely/Somewhat Disagree

**Patients are concerned when PII is shared**  
**Solution: disclosure control can help**

## Figure 2. Data Sharing Comfort (n=3516): Local Privacy



I am comfortable with researchers not directly involved in my care accessing my electronic health data for research purposes



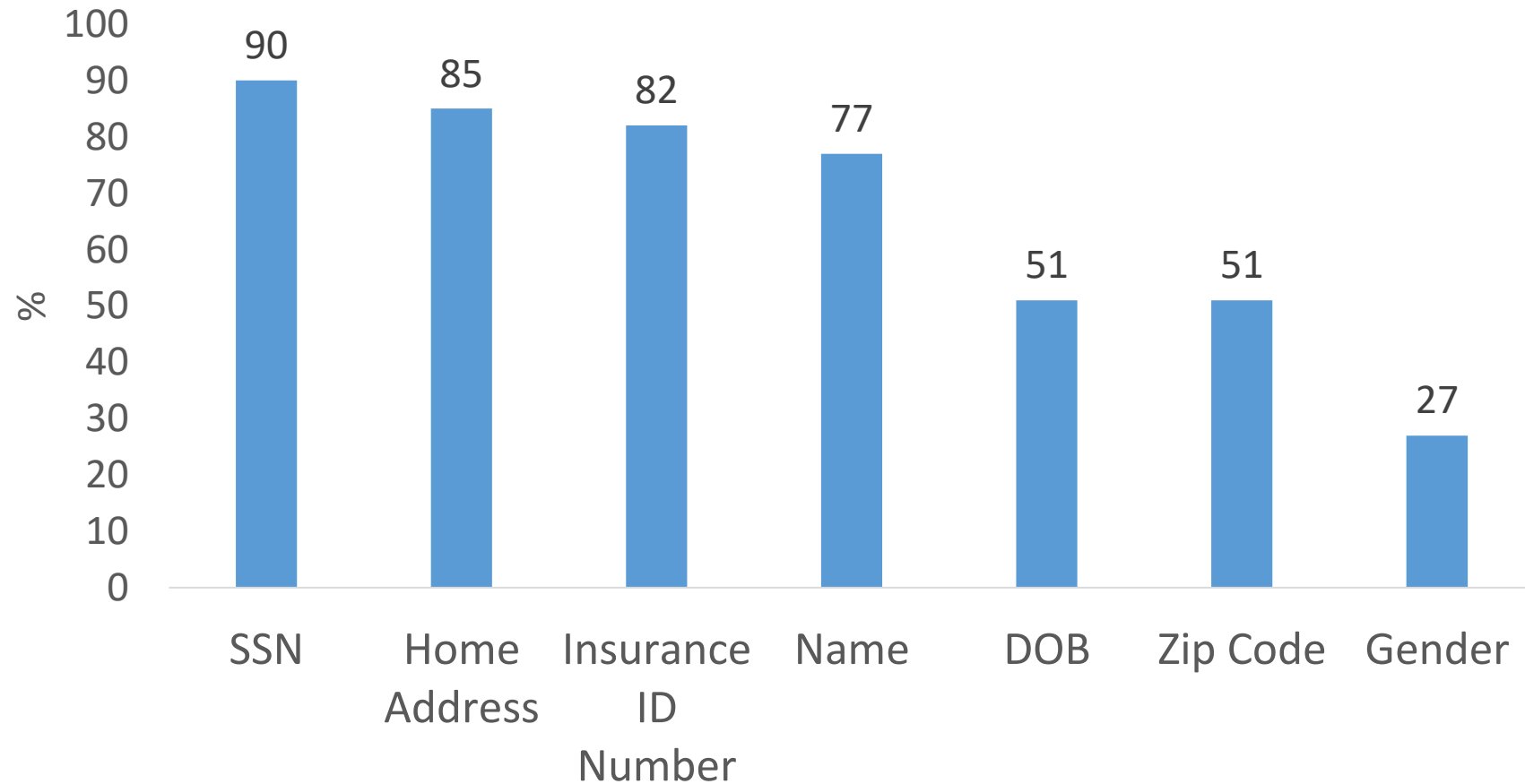
I am comfortable with someone I know (e.g., friend, neighbor, coworker) who is a researcher accessing my electronic health data for research purposes

■ Completely/Somewhat Agree ■ Moderately Agree ■ Completely/Somewhat Disagree

More patients are concerned when someone who can recognize them (e.g., someone who knows them) accesses EHR for research

Solution: disclosure control can reduce people who know you, recognizing you

Figure 3. % reporting they would be “extremely” or “much more comfortable” with removal of the following identifiers (n=3516):



**What attributes are they most concerned about ?**

**Solution: Focus on masking IDs, names, and addresses**



# Goal: Build consensus on template IRB application & DUA when using PPIRL framework

- The questions we plan to ask at the NGT session are:
  1. What do you perceive as the benefits when using the PPIRL framework for record linkage?
    - Potentially, allows for linking data that would otherwise not be possible.
    - Encourages use of only needed information, minimizing risk
    - Minimizes risk of re-identification, reduces risk of breach of confidentiality
  2. What do you perceive as the risks when using the PPIRL framework for record linkage?
    - Mislinking risks
    - size and quality of data matters; a bad database makes linking difficult when data is masked
    - Disproportionate data sampling, which leads do an increase in bias
  3. What other information would you like or need to know when reviewing the IRB application for research?
  4. When using the PPIRL framework, what information is needed in the DUA?
    - How to communicate risk to lawyers, so that the risk is stated in the DUA accurately
    - Since the DUA is fixed, how can the software adjust to the DUA
    - Expert determination