



ai4se.net

How to Read Less:

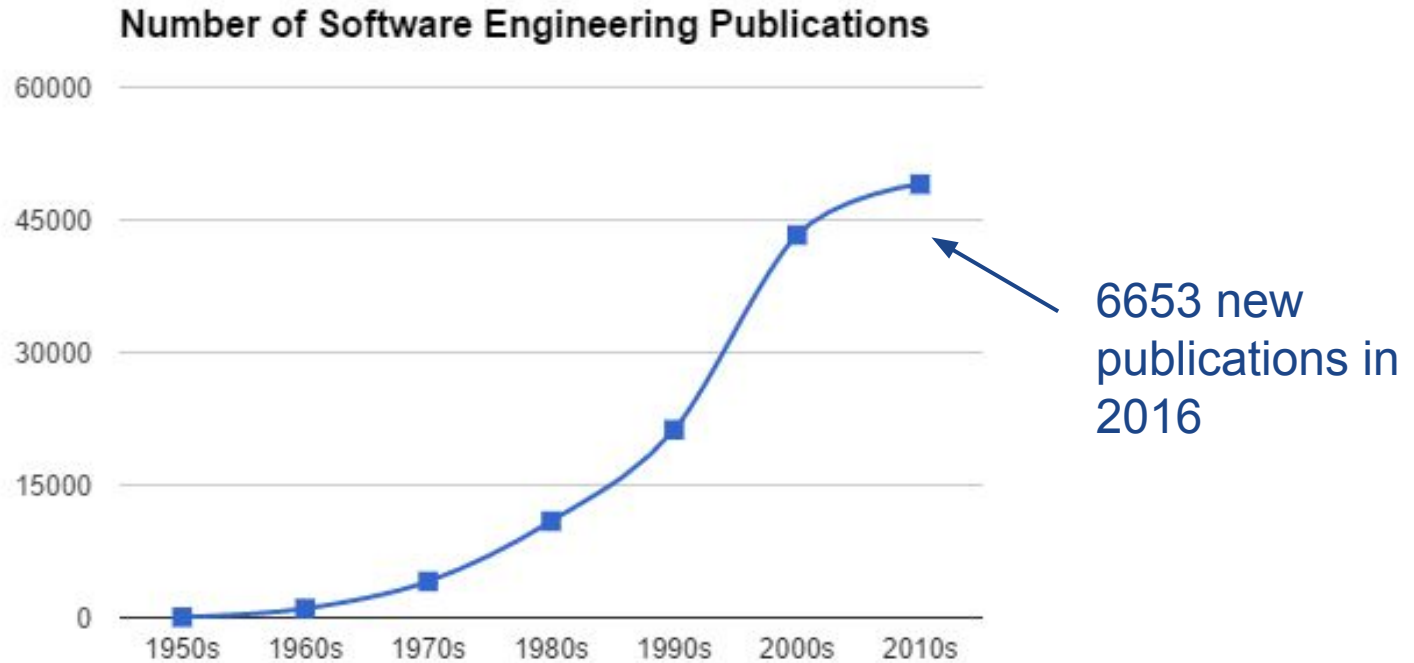
On the Benefit of Human-in-the-loop Incremental Learning for Systematic Literature Reviews

Zhe Yu
zyu9@ncsu.edu

Ph.D. Written Preliminary Exam
April 2017

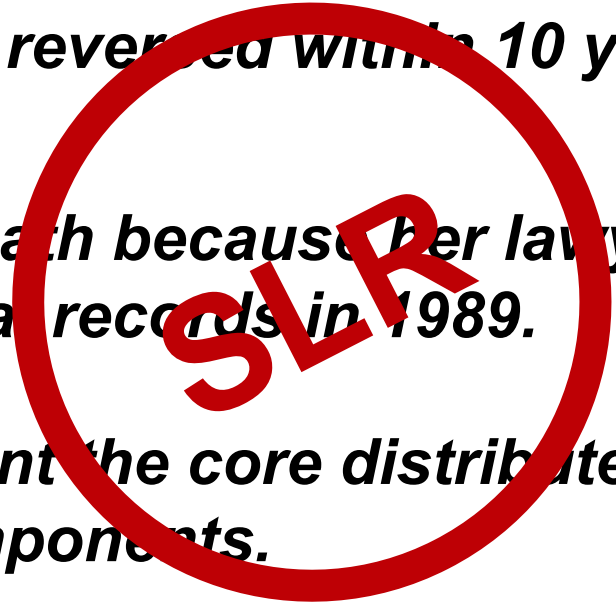
Committee:
Dr. Tim Menzies (advisor)
Dr. Xipeng Shen
Dr. William Enck

Too Much to Read



Not Reading Enough

- *146 medical practises were reversed within 10 years.*
- *A lady was sentenced to death because her lawyer failed to find related hospital records in 1989.*
- *Software developers reinvent the core distributed software concepts and components.*



Done

- *Apply human-in-the-loop method to facilitate SLRs*
- *FASTREAD, a better method than state-of-the-art*
- *Save 90% cost to retrieve 90% relevant studies*
- *A tool to implement FASTREAD*

OpenSource

Paper:

- Submitted to EMSE
- [arxiv](#)

Tool:

- [SeaCraft Zenodo](#)
- [Github](#)

The screenshot shows a web application interface for document classification. The interface is divided into several sections:

- File Selection (A):** A 'Choose File' button is highlighted with a red arrow. Below it are 'Next', 'Auto Review', and 'Random' buttons.
- Document List (B):** A list of documents is displayed, with the sixth item, 'Estimating software fault-proneness for tuning testing activities', highlighted in yellow and marked with a red arrow (B).
- Classification Control (C, D):** Radio buttons for 'Relevant', 'Irrelevant', and 'Undetermined' are shown, along with a 'Submit' button (D). The 'True Label: yes' is displayed below.
- Plot (E, F, G, H, I, J):** A line graph titled 'Documents Coded: 92/280 (8911)' shows 'Relevant Found' on the y-axis (0 to 100) and 'Documents Reviewed' on the x-axis (0 to 300). The plot shows a curve that rises from 0 to approximately 90. The plot area is marked with red arrows (E, F, G, H, I, J) pointing to various points on the curve.

Todo

- *More data*
- *More detailed, real data*
- *Massive data*

Outline

- *Background*
- *Method*
- *Experiment*
- *Result*
- *Conclusion and Future works*



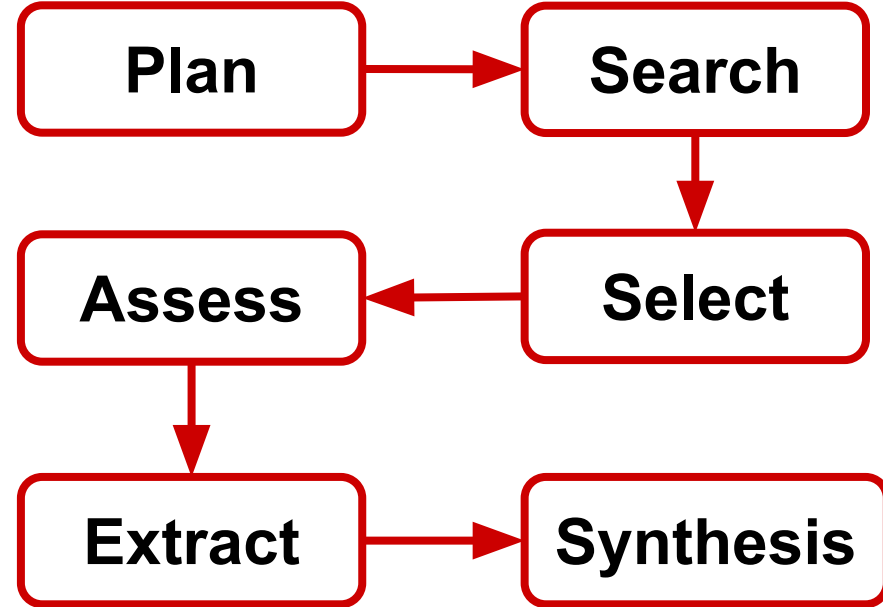
Outline

- ***Background***
 - Systematic Literature Review
 - Primary Study Selection (Select)
- ***Method***
- ***Experiment***
- ***Result***
- ***Conclusion and Future works***



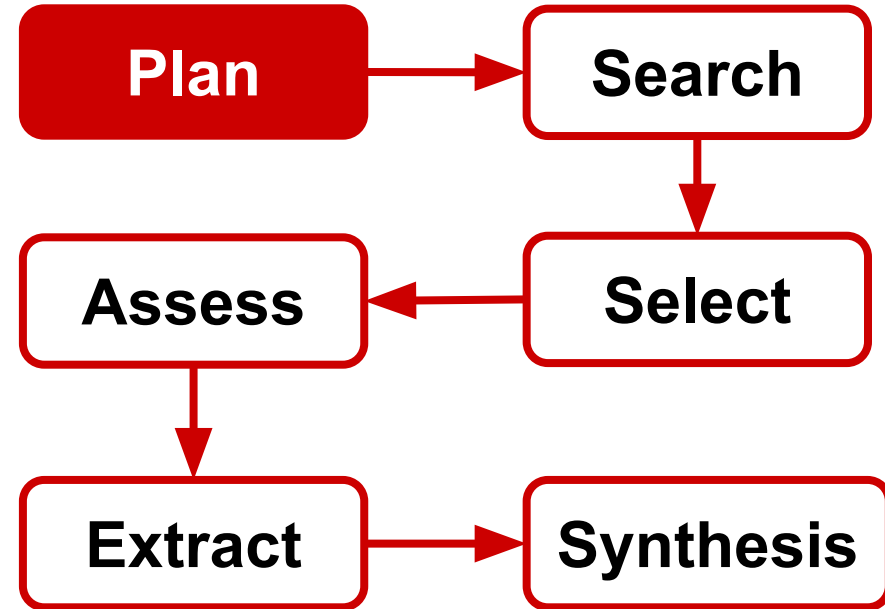
Systematic Literature Review (SLR)

SLR: A systematic guide to review literature



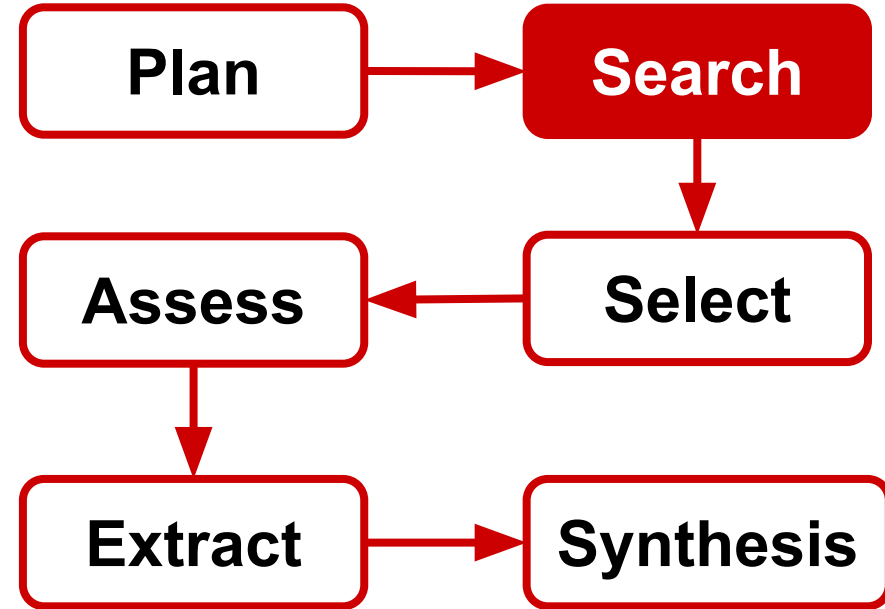
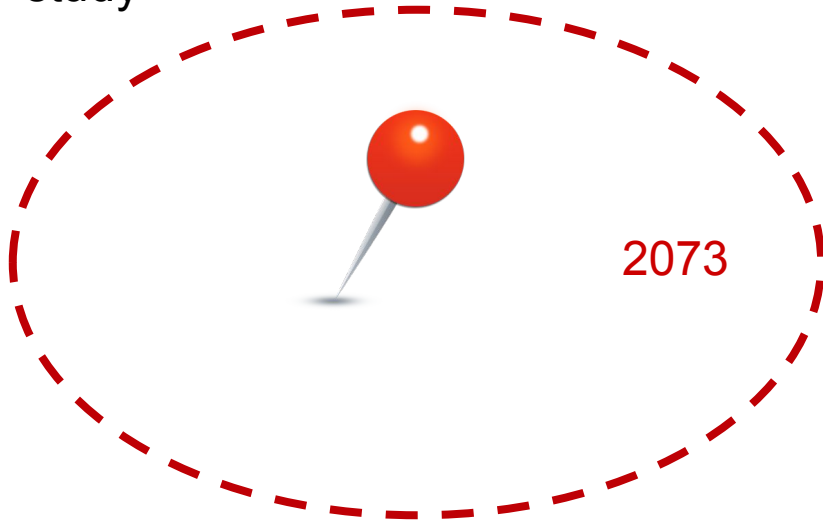
Defect Prediction Papers, from
2000 to 2010 [Hall'12]

- Context
- Feature
- Modeling



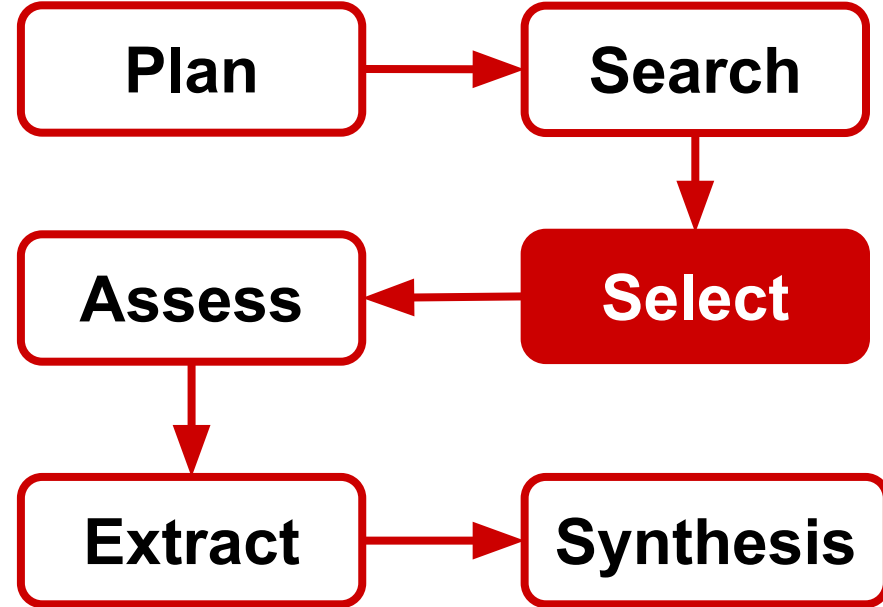
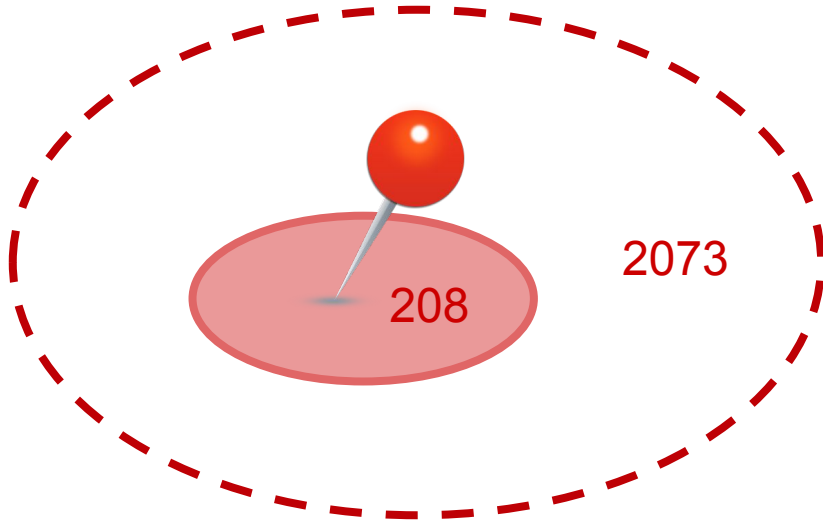
Search:

(Fault* OR bug* OR defect* OR errors
OR corrections OR corrective OR fix*) in
title only AND (Software) anywhere in
study



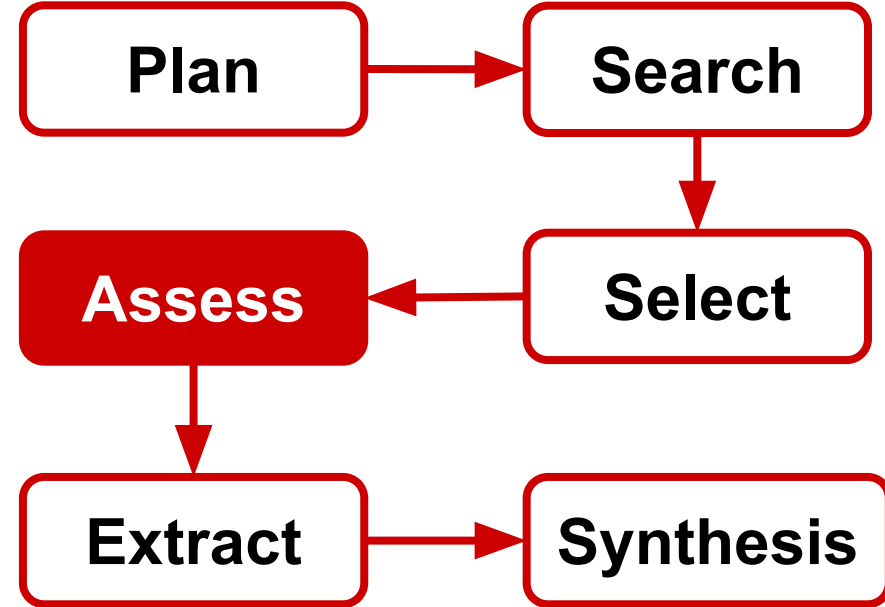
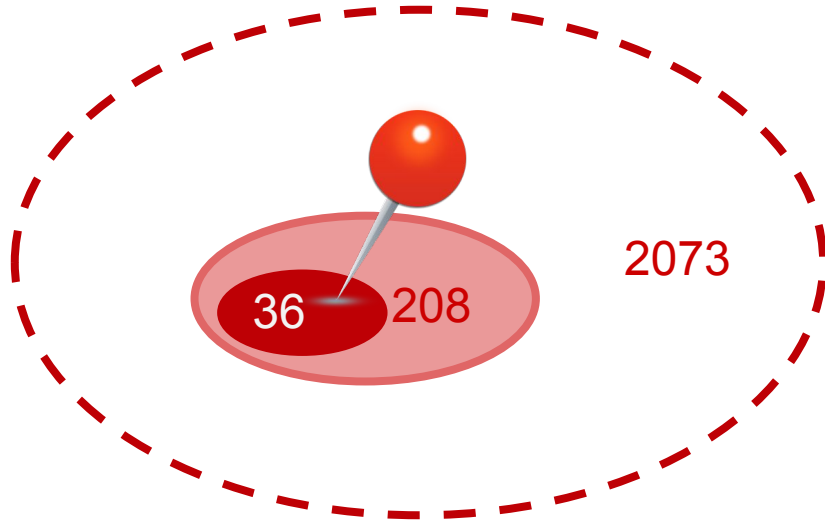
Exclude:

- not an empirical study
- not in software systems
- main output is not faults in code



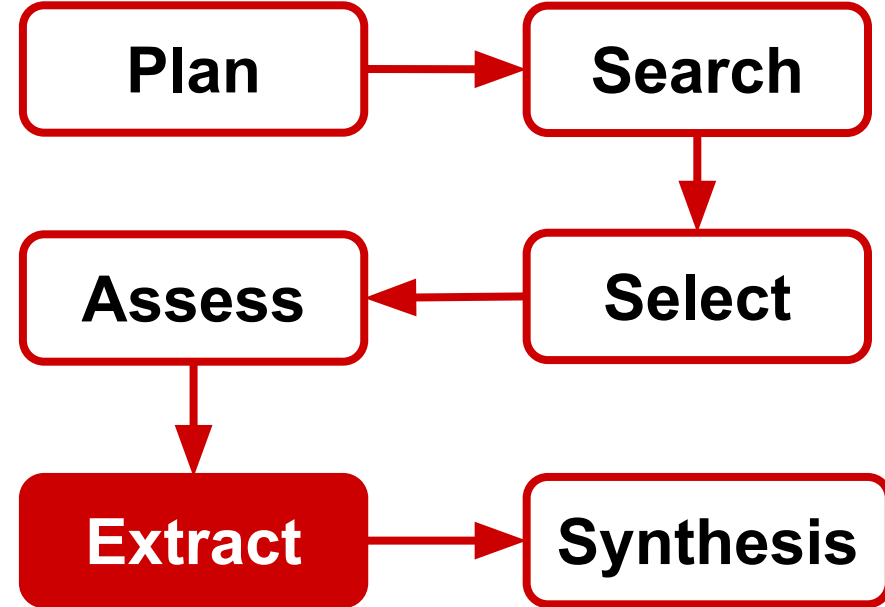
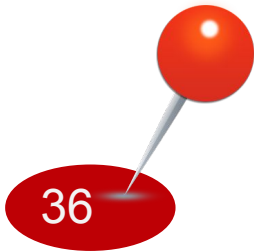
Assess:

- is a prediction study
- sufficient contextual information
- sufficient model building information
- has data



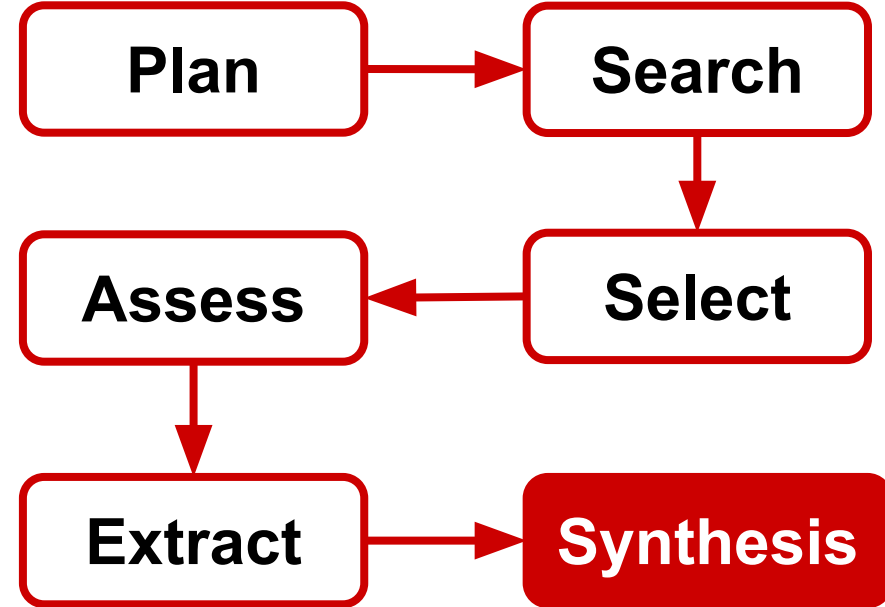
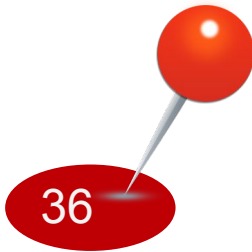
Extract:

- Context data
- Qualitative data
- Quantitative data

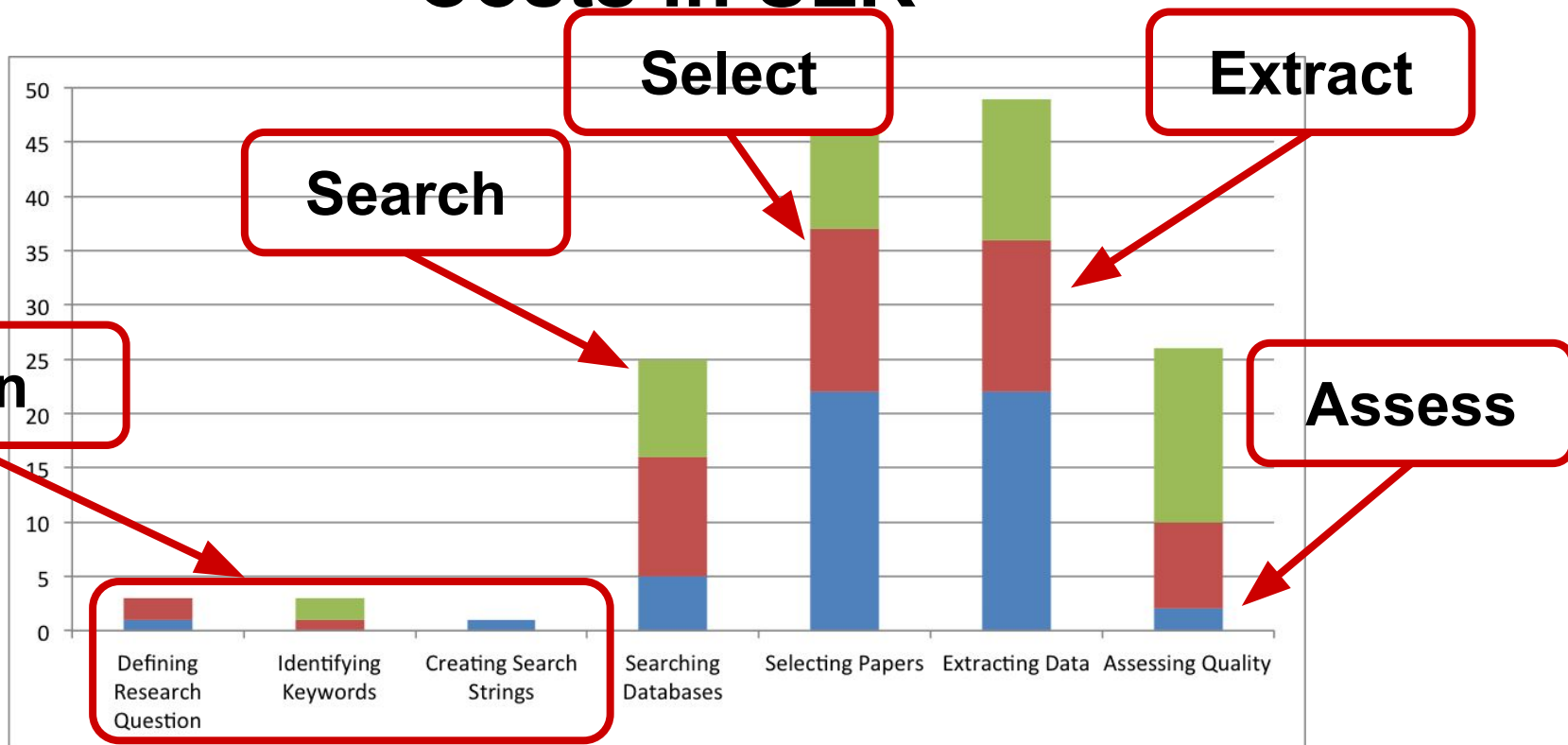


Synthesis:

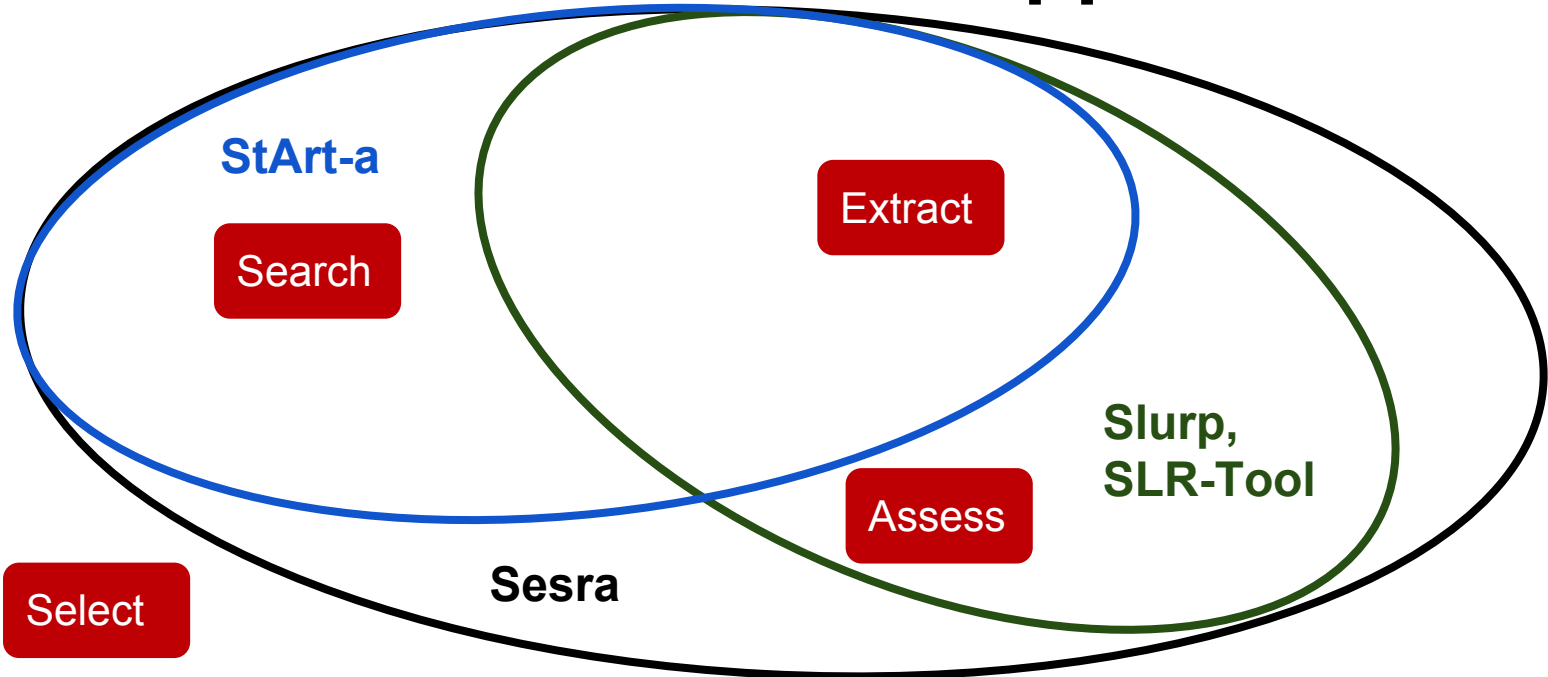
- **What features are used most frequently?** LOC, OO, etc.
- **Which model works best?** Naive Bayes, Decision Tree, Logistic Regression, etc.



Costs in SLR



Automated Tool Support



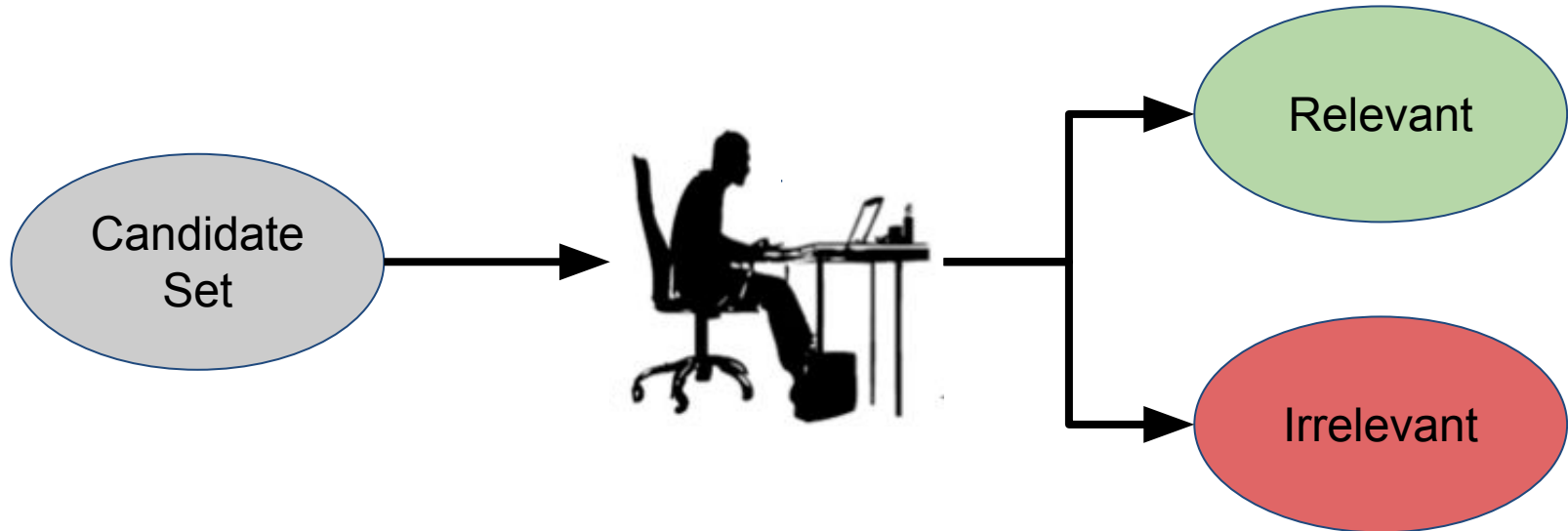
Human-in-the-loop, VTM, Snowballing

Outline

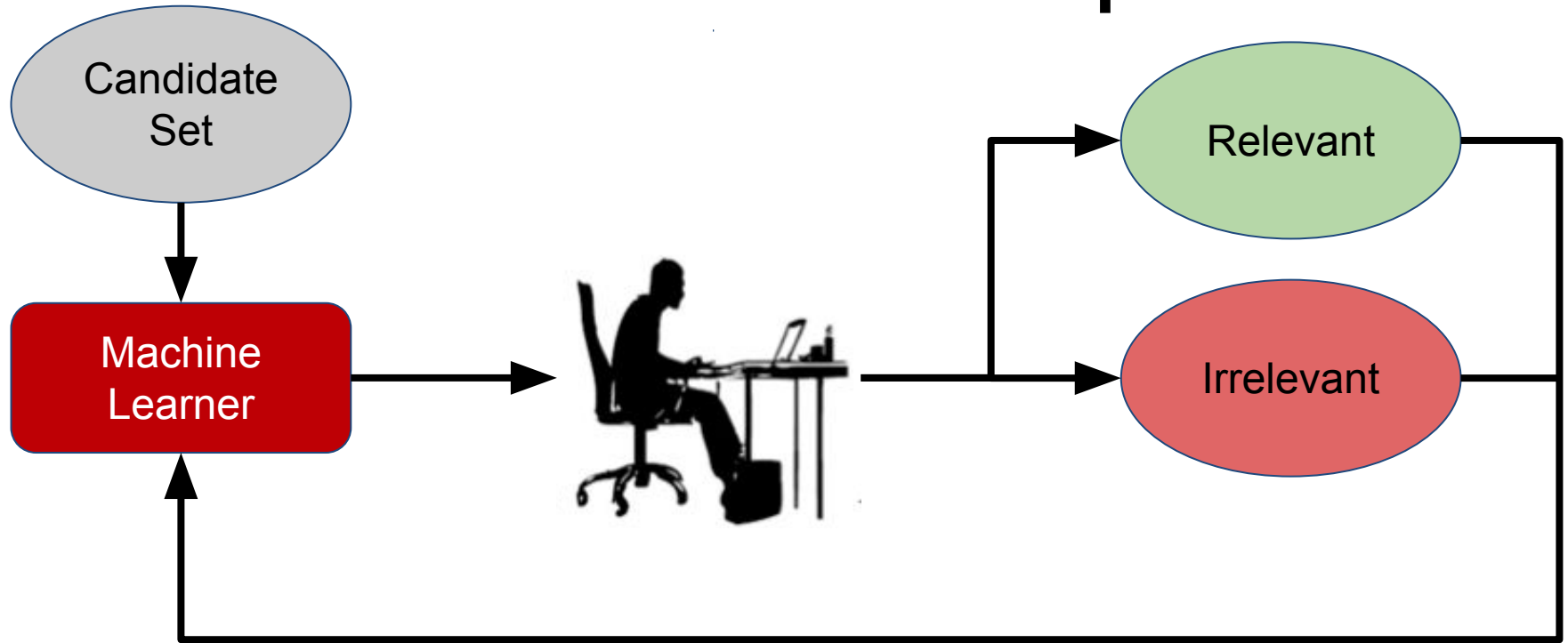
- ***Background***
 - Systematic Literature Review
 - Primary Study Selection (Select)
- ***Method***
- ***Experiment***
- ***Result***
- ***Conclusion and Future works***



Select: Linear Review



Human-in-the-loop

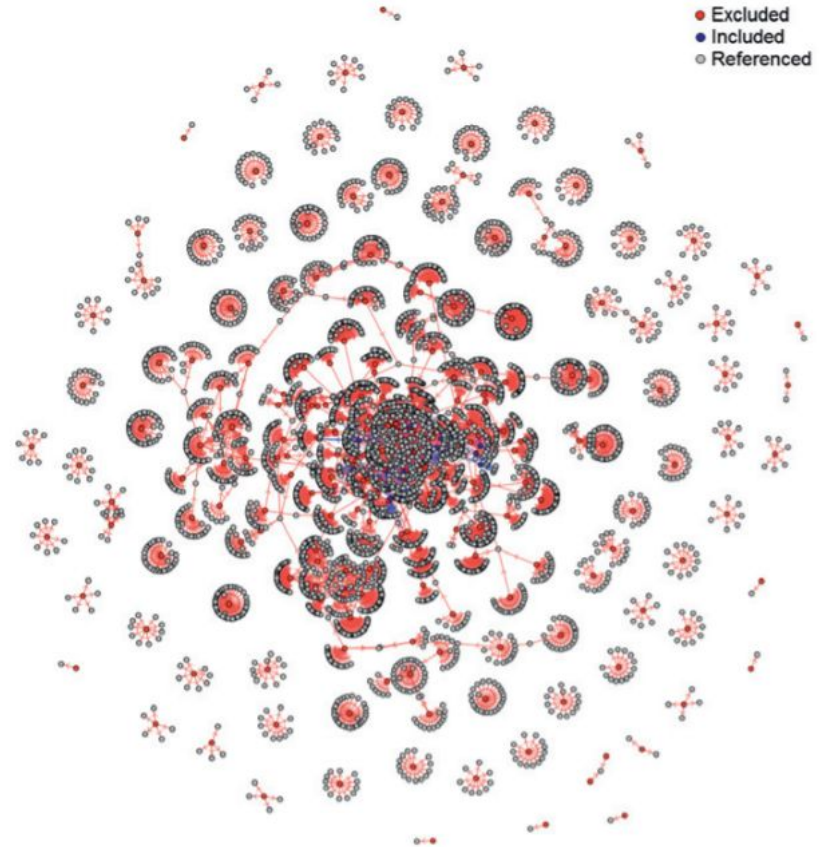
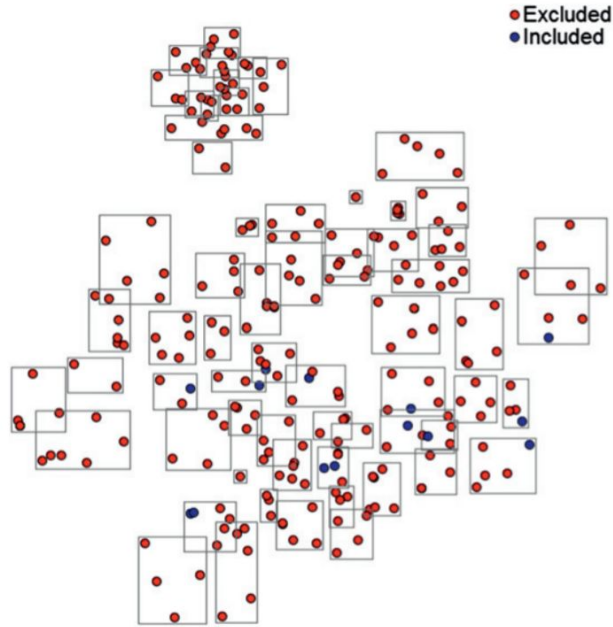


Sheng, Provost, Ipeirotis. 2008. "Get another label? improving data quality and data mining using multiple, noisy labelers", KDD '08

Wallace, Byron C., et al. "Semi-automated screening of biomedical citations for systematic reviews." BMC bioinformatics 11.1 (2010): 55.

Cormack, Gordon V., and Maura R. Grossman. "Evaluation of machine-learning protocols for technology-assisted review in electronic discovery." SIGIR'14.

Visual Text Mining



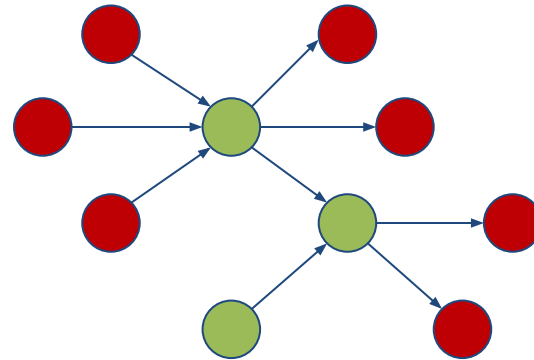
Snowballing

Forward:

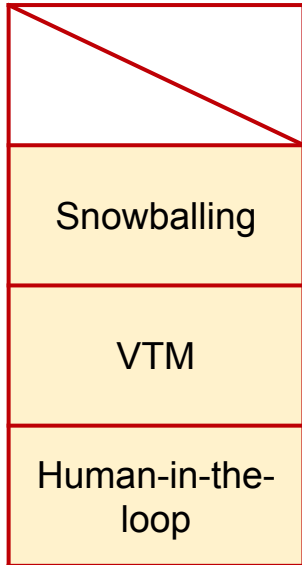
Papers cite the known one.

Backward:

The known one's references



Techniques to Select



Wohlin, Claes. "Guidelines for snowballing in systematic literature studies and a replication in software engineering." EASE'14.

Felizardo, Katia R., et al. "visual analysis approach to validate the selection review of primary studies" Information and Software Technology 54.10 (2012): 1079-1091.

Wallace, Byron C., et al. "Semi-automated screening of biomedical citations for systematic reviews." BMC bioinformatics 11.1 (2010): 55.

Cormack, Gordon V., and Maura R. Grossman. "Evaluation of machine-learning protocols for technology-assisted review in electronic discovery." SIGIR'14.

Techniques to Select

	Machine Learning?
Snowballing	no
VTM	Unsupervised Learning
Human-in-the-loop	Active Learning

Wohlin, Claes. "Guidelines for snowballing in systematic literature studies and a replication in software engineering." EASE'14.

Felizardo, Katia R., et al. "visual analysis approach to validate the selection review of primary studies" Information and Software Technology 54.10 (2012): 1079-1091.

Wallace, Byron C., et al. "Semi-automated screening of biomedical citations for systematic reviews." BMC bioinformatics 11.1 (2010): 55.

Cormack, Gordon V., and Maura R. Grossman. "Evaluation of machine-learning protocols for technology-assisted review in electronic discovery." SIGIR'14.

Techniques to Select

	Machine Learning?	Require initial papers?
Snowballing	no	yes
VTM	Unsupervised Learning	no
Human-in-the-loop	Active Learning	no

Wohlin, Claes. "Guidelines for snowballing in systematic literature studies and a replication in software engineering." EASE'14.

Felizardo, Katia R., et al. "visual analysis approach to validate the selection review of primary studies" Information and Software Technology 54.10 (2012): 1079-1091.

Wallace, Byron C., et al. "Semi-automated screening of biomedical citations for systematic reviews." BMC bioinformatics 11.1 (2010): 55.

Cormack, Gordon V., and Maura R. Grossman. "Evaluation of machine-learning protocols for technology-assisted review in electronic discovery." SIGIR'14.

Techniques to Select

	Machine Learning?	Require initial papers?	Easy to Validate?
Snowballing	no	yes	yes
VTM	Unsupervised Learning	no	no
Human-in-the-loop	Active Learning	no	yes

Wohlin, Claes. "Guidelines for snowballing in systematic literature studies and a replication in software engineering." EASE'14.

Felizardo, Katia R., et al. "visual analysis approach to validate the selection review of primary studies" Information and Software Technology 54.10 (2012): 1079-1091.

Wallace, Byron C., et al. "Semi-automated screening of biomedical citations for systematic reviews." BMC bioinformatics 11.1 (2010): 55.

Cormack, Gordon V., and Maura R. Grossman. "Evaluation of machine-learning protocols for technology-assisted review in electronic discovery." SIGIR'14.

Techniques to Select

	Machine Learning?	Require initial papers?	Easy to Validate?	Scale up?
Snowballing	no	yes	yes	?
VTM	Unsupervised Learning	no	no	?
Human-in-the-loop	Active Learning	no	yes	yes

Wohlin, Claes. "Guidelines for snowballing in systematic literature studies and a replication in software engineering." EASE'14.

Felizardo, Katia R., et al. "visual analysis approach to validate the selection review of primary studies" Information and Software Technology 54.10 (2012): 1079-1091.

Wallace, Byron C., et al. "Semi-automated screening of biomedical citations for systematic reviews." BMC bioinformatics 11.1 (2010): 55.

Cormack, Gordon V., and Maura R. Grossman. "Evaluation of machine-learning protocols for technology-assisted review in electronic discovery." SIGIR'14.

Techniques to Select

	Machine Learning?	Require initial papers?	Easy to Validate?	Scale up?	Tested on SE data?
Snowballing	no	yes	yes	?	yes
VTM	Unsupervised Learning	no	no	?	yes
Human-in-the-loop	Active Learning	no	yes	yes	no

Wohlin, Claes. "Guidelines for snowballing in systematic literature studies and a replication in software engineering." EASE'14.

Felizardo, Katia R., et al. "visual analysis approach to validate the selection review of primary studies" Information and Software Technology 54.10 (2012): 1079-1091.

Wallace, Byron C., et al. "Semi-automated screening of biomedical citations for systematic reviews." BMC bioinformatics 11.1 (2010): 55.

Cormack, Gordon V., and Maura R. Grossman. "Evaluation of machine-learning protocols for technology-assisted review in electronic discovery." SIGIR'14.

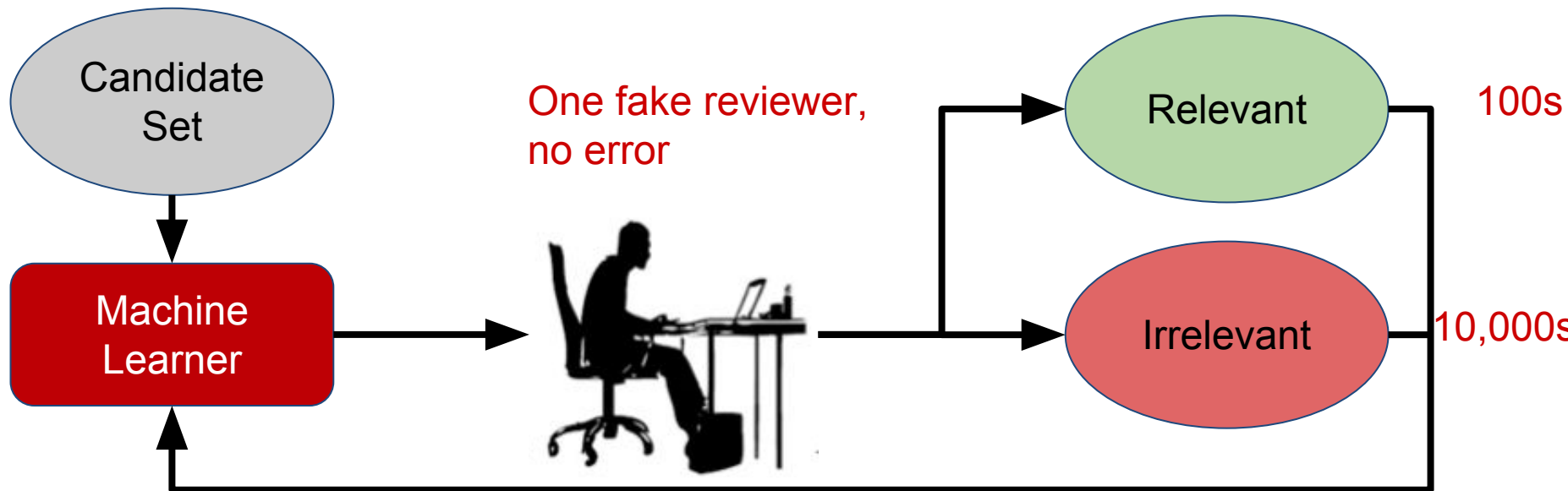
Outline

- *Background*
- ***Method***
 - Problem Statement
 - State-of-the-art Methods
 - Algorithm Code
- *Experiment*
- *Result*
- *Conclusion and Future works*



Problem Statement

Title+Abstract,
unlabeled, 10,000s



Wallace, Byron C., et al. "Semi-automated screening of biomedical citations for systematic reviews." BMC bioinformatics 11.1 (2010): 55.

Cormack, Gordon V., and Maura R. Grossman. "Evaluation of machine-learning protocols for technology-assisted review in electronic discovery." SIGIR'14.

Hall, Tracy, et al. "A systematic literature review on fault prediction performance in software engineering." TSE'12.

Outline

- *Background*
- **Method**
 - Problem Statement
 - State-of-the-art Methods
 - Algorithm Code
- *Experiment*
- *Result*
- *Conclusion and Future works*



State-of-the-art

Medicine:

[Wallace'10] *Core algorithm stays unchanged in subsequent works.*

Legal:

[Cormack'14] *Still state-of-the-art in legal domain.*

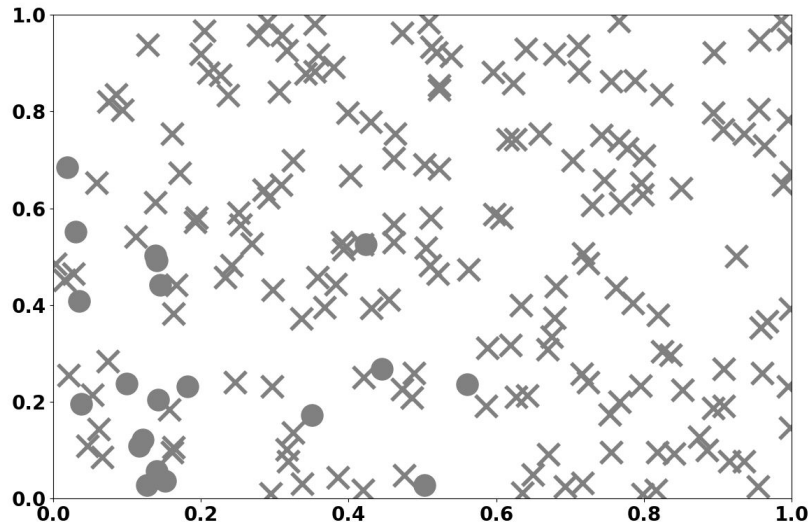
Outline

- *Motivation and Background*
- **Method**
 - Problem Statement
 - State-of-the-art
 - Algorithm Code
- *Experiment*
- *Result*
- *Conclusion and Future works*

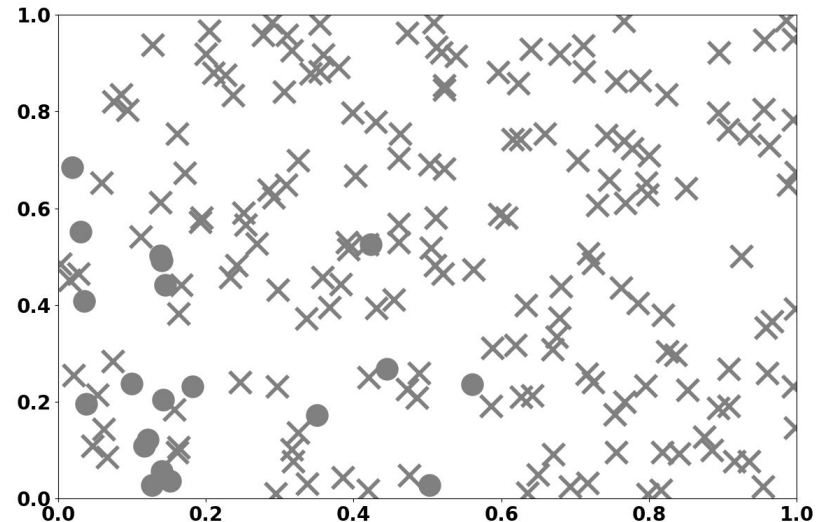


Initial Candidates

Wallace'10 (0/0)

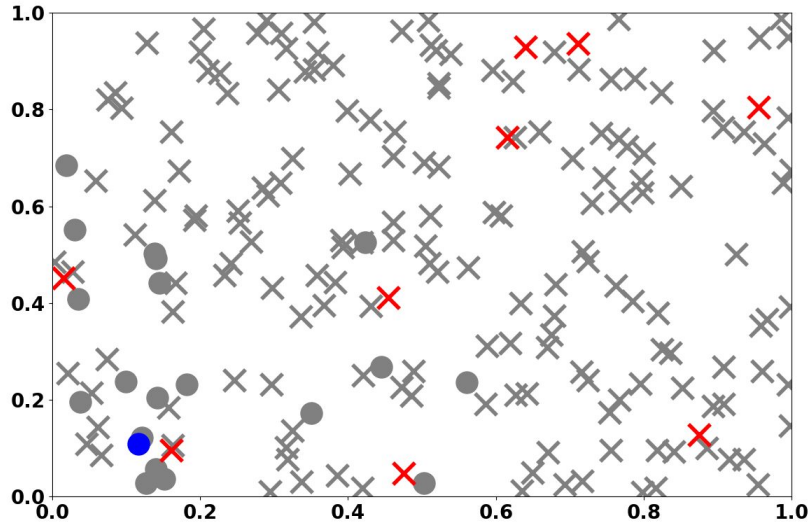


Cormack'14 (0/0)

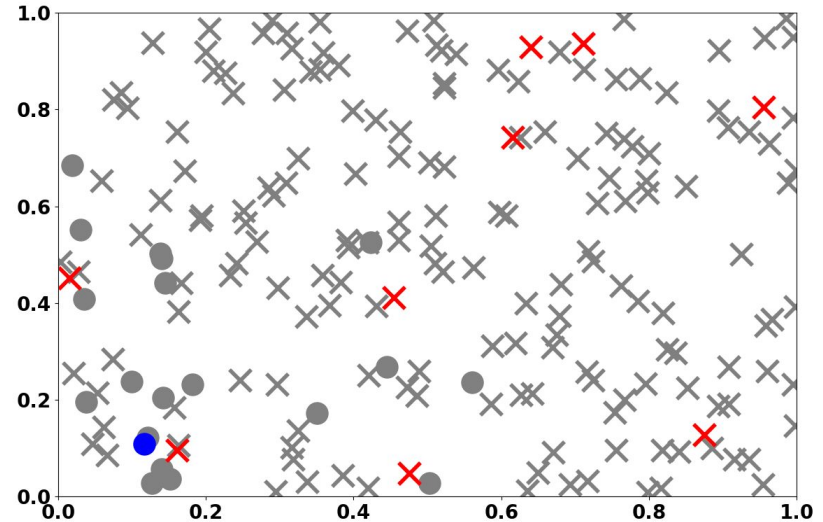


Random Sampling

Wallace'10 (1/10)

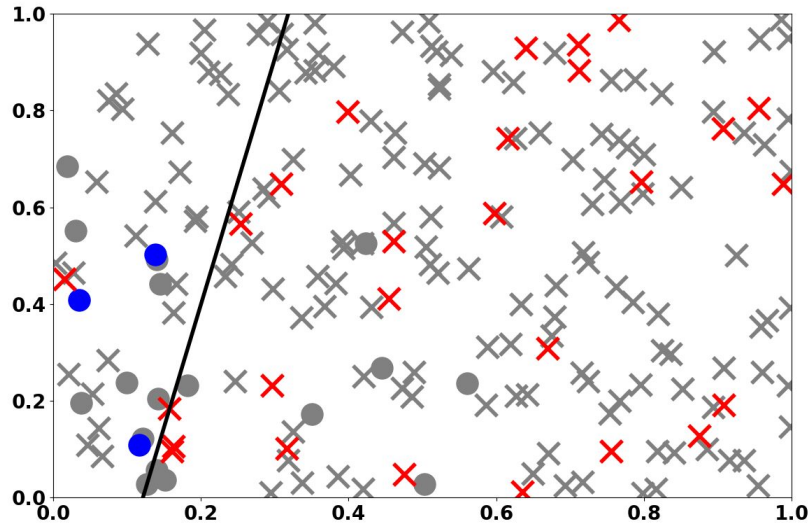


Cormack'14 (1/10)

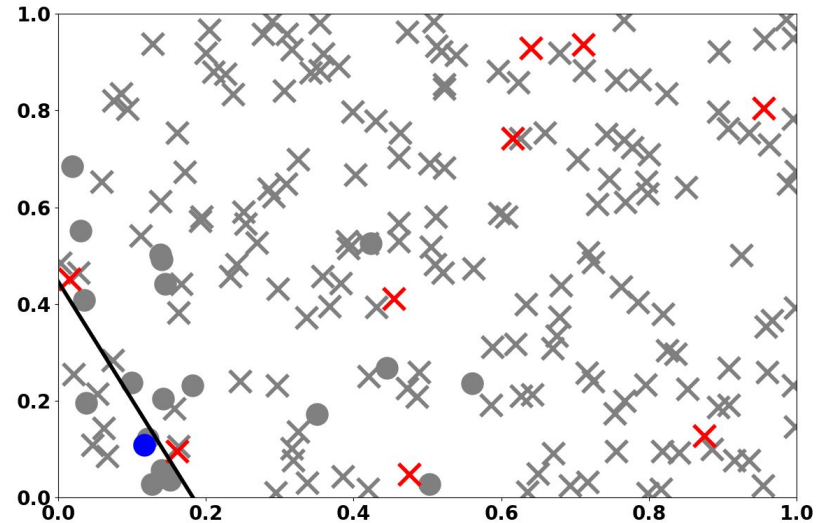


Start Point

Wallace'10 (3/30)

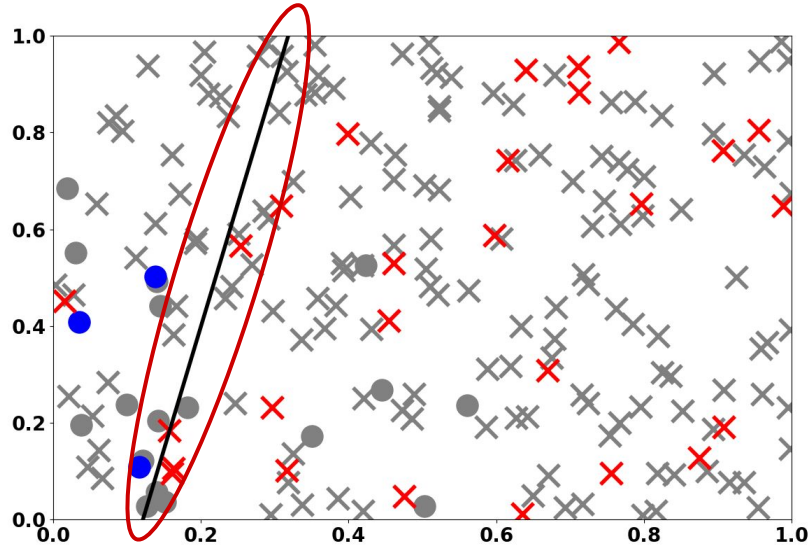


Cormack'14 (1/10)



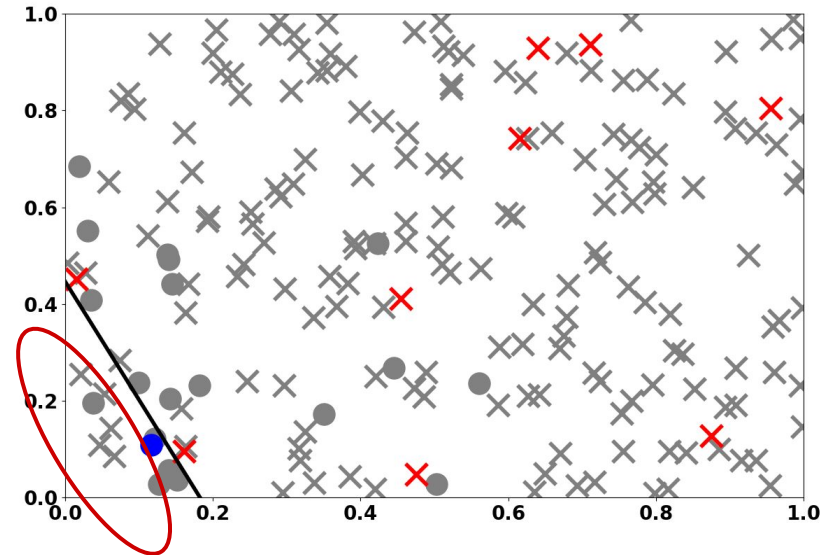
Query Strategy

Wallace'10 (3/30)



Uncertainty Sampling

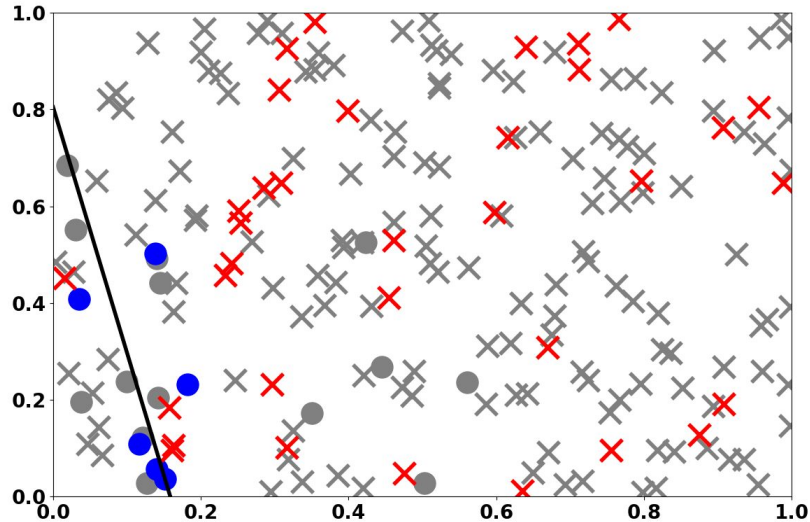
Cormack'14 (1/10)



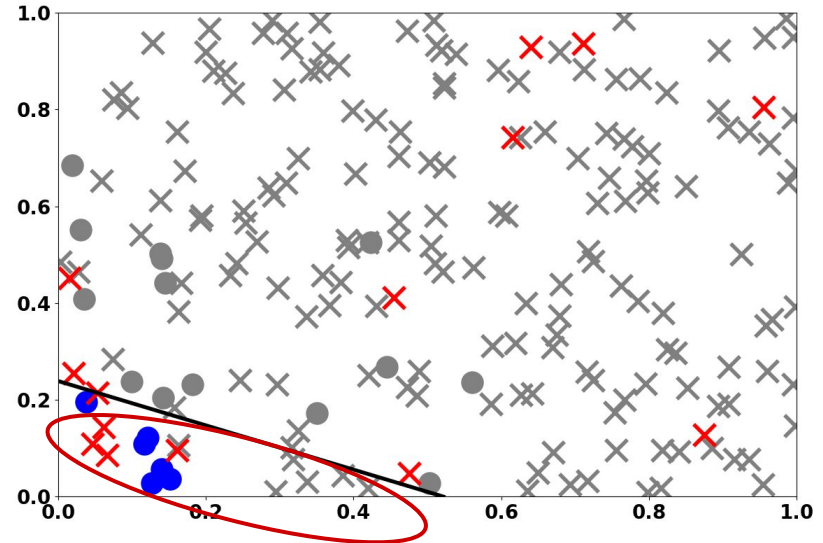
Certainty Sampling

Data Balancing

Wallace'10 (6/40)



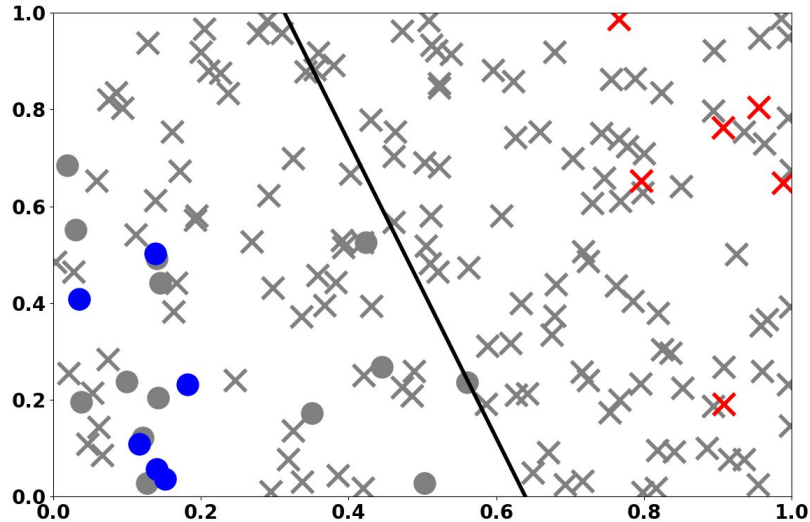
Cormack'14 (6/20)



Certainty Sampling

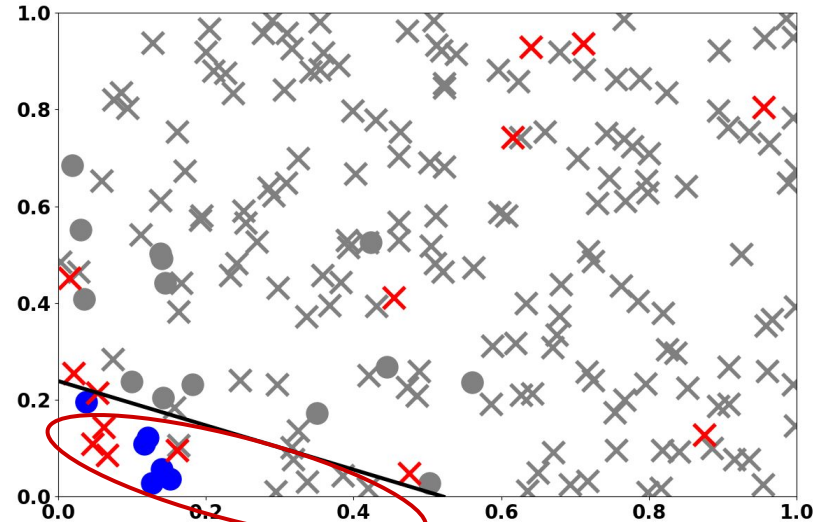
Data Balancing

Wallace'10 (6/40)



Aggressive Undersampling

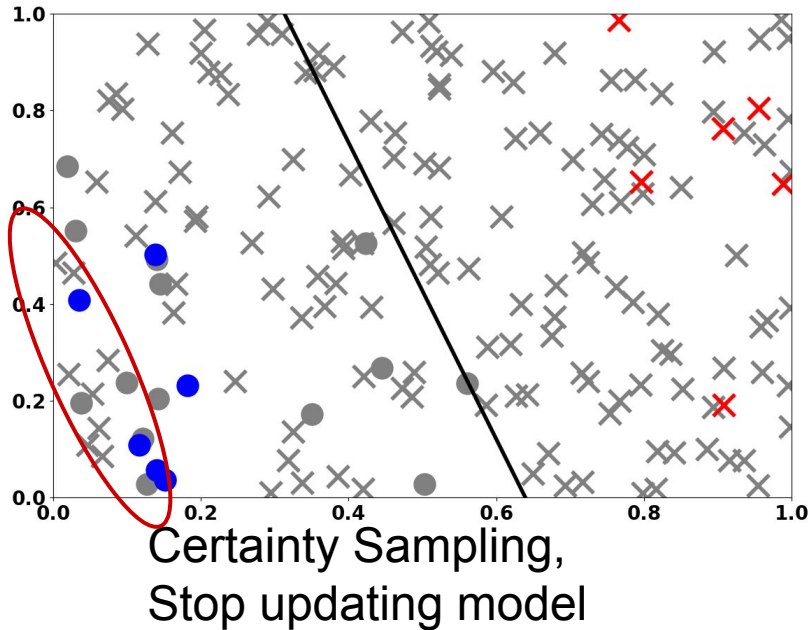
Cormack'14 (6/20)



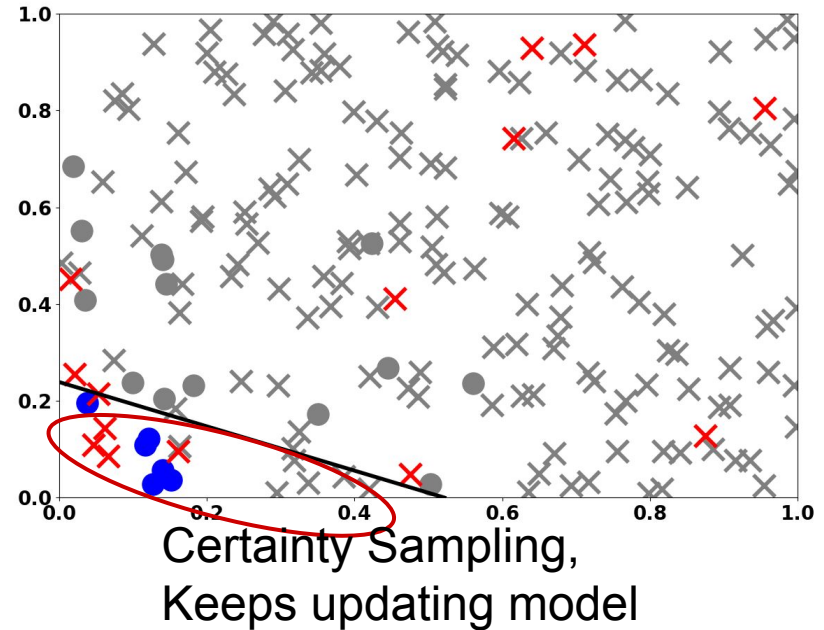
Certainty Sampling

Continuity

Wallace'10 (6/40)



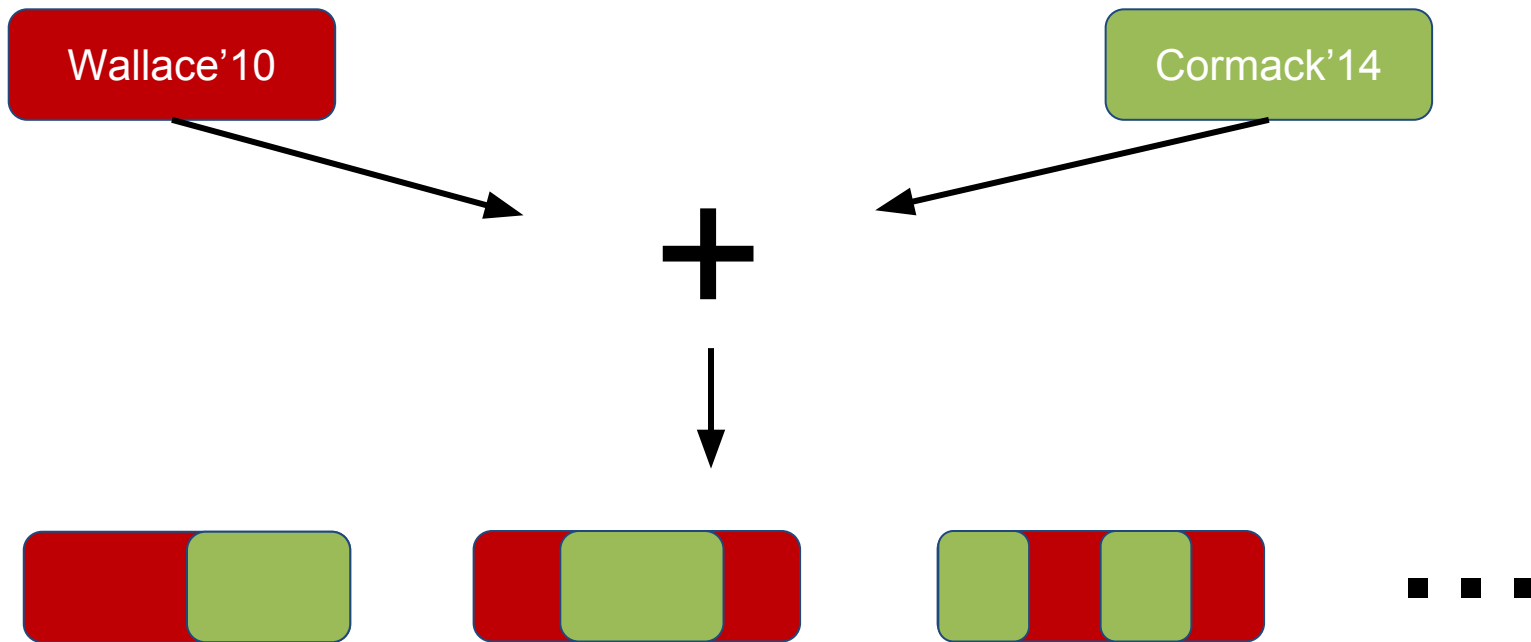
Cormack'14 (6/20)



Algorithm Code

	Start point	Query Strategy	Continuity	Data Balancing
Wallace'10	late (P)	uncertainty sampling (U)	stop training (S)	aggressive undersampling (A)
Cormack'14	early (\bar{P})	certainty sampling (\bar{U})	non-stop (\bar{S})	none (\bar{A})

Refactor



Outline

- *Background*
- *Method*
- ***Experiment***
 - Data
 - Framework
 - Evaluation
- *Result*
- *Conclusion and Future works*



Data Sets

Wahono'15

7002

62

Hall'12

8911

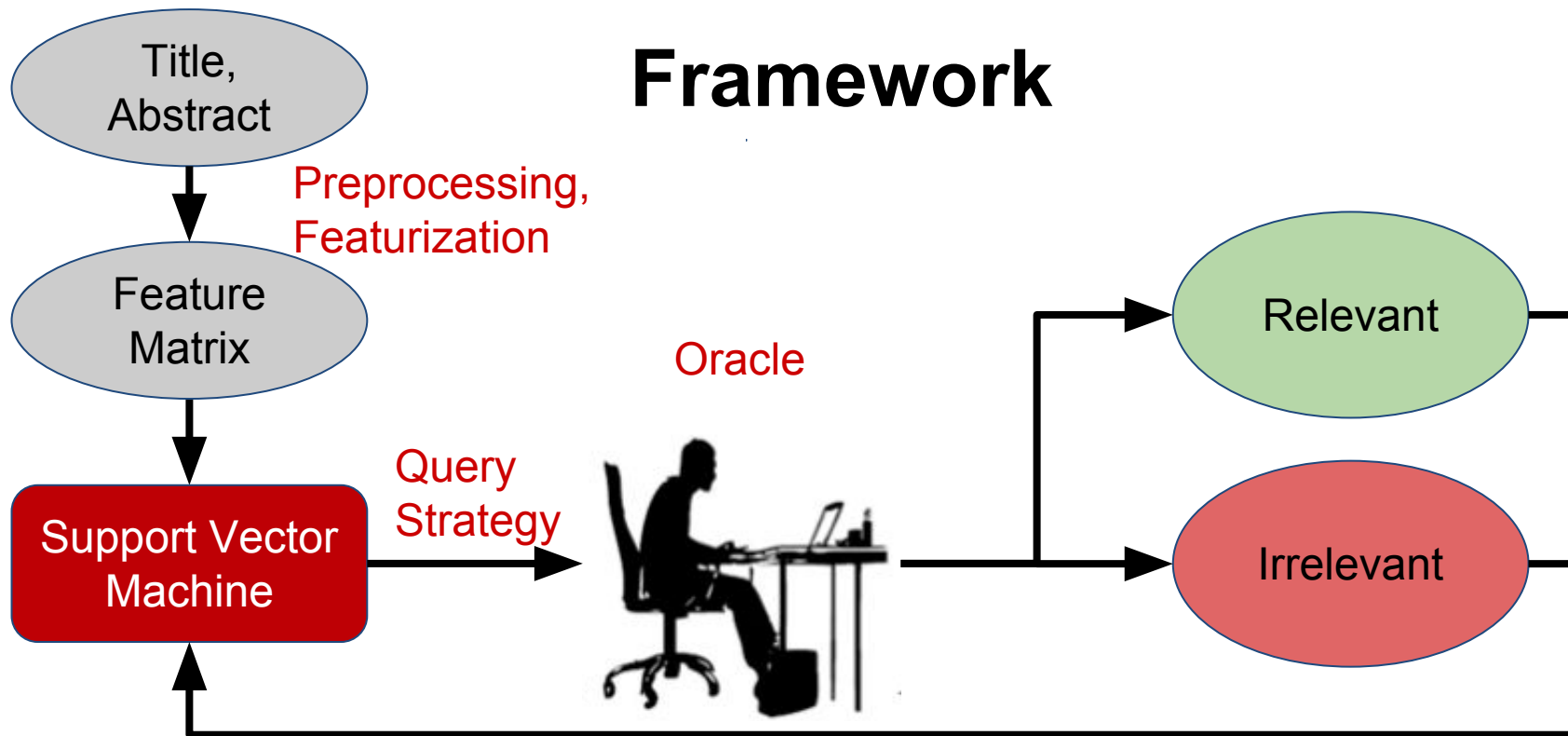
106

Outline

- *Background*
- *Method*
- ***Experiment***
 - Data
 - Framework
 - Evaluation
- *Result*
- *Conclusion and Future works*



Framework

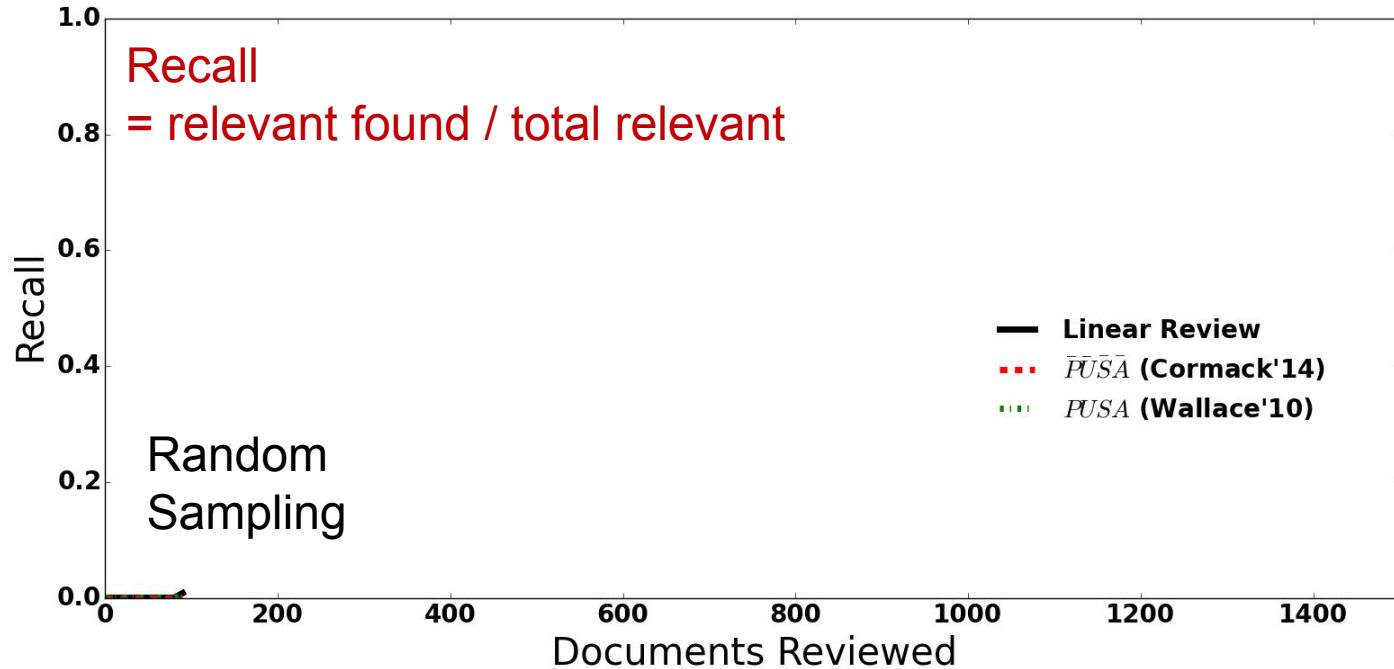


Outline

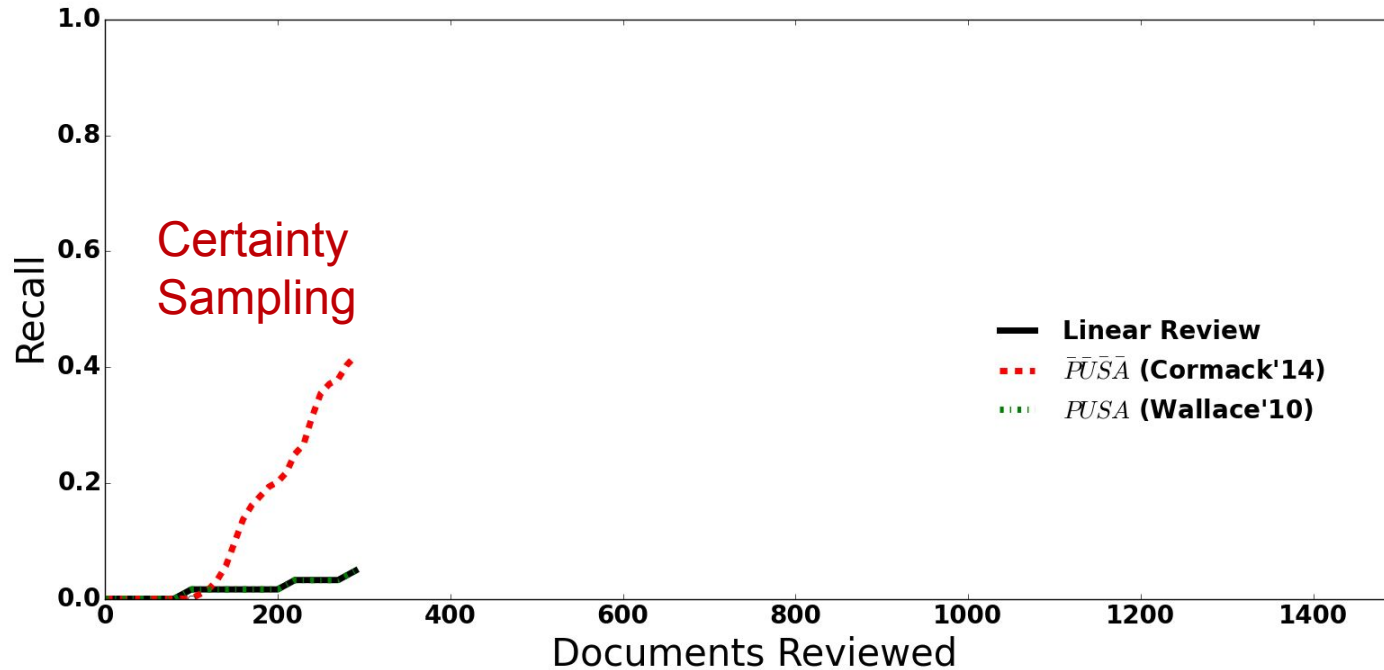
- *Background*
- *Method*
- ***Experiment***
 - Data
 - Framework
 - Evaluation
- *Result*
- *Conclusion and Future works*



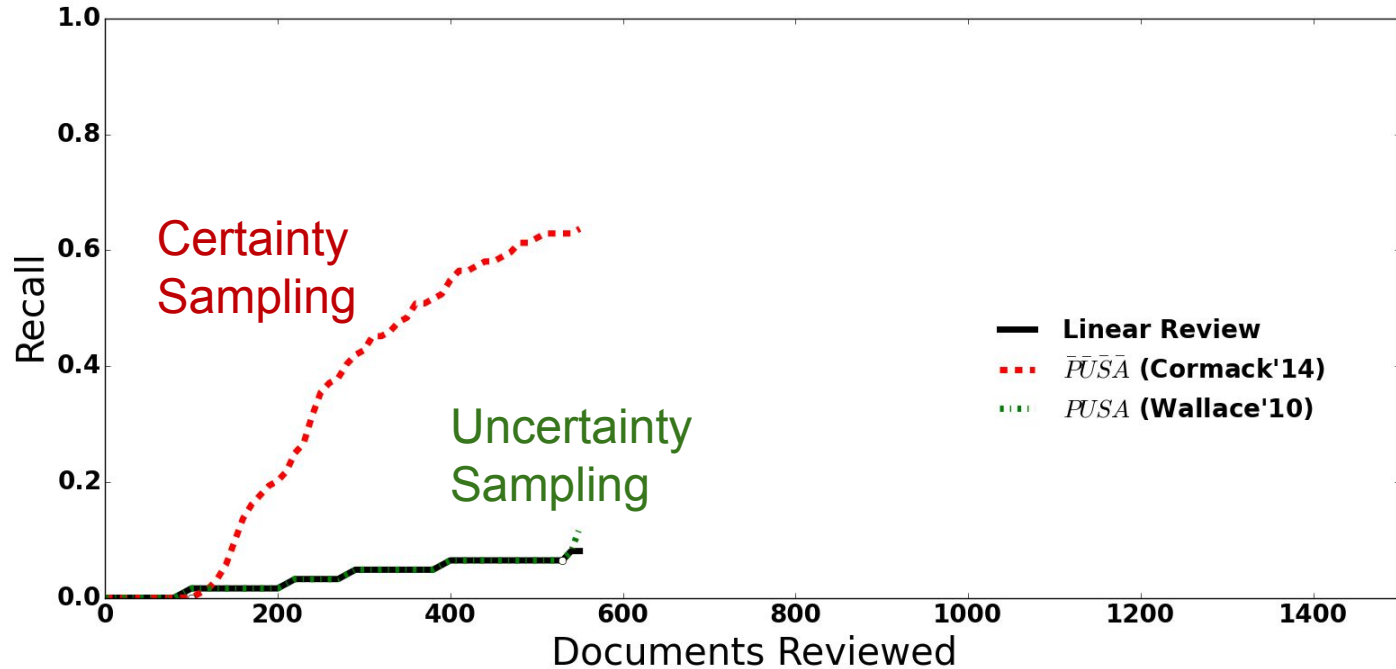
Evaluation 1



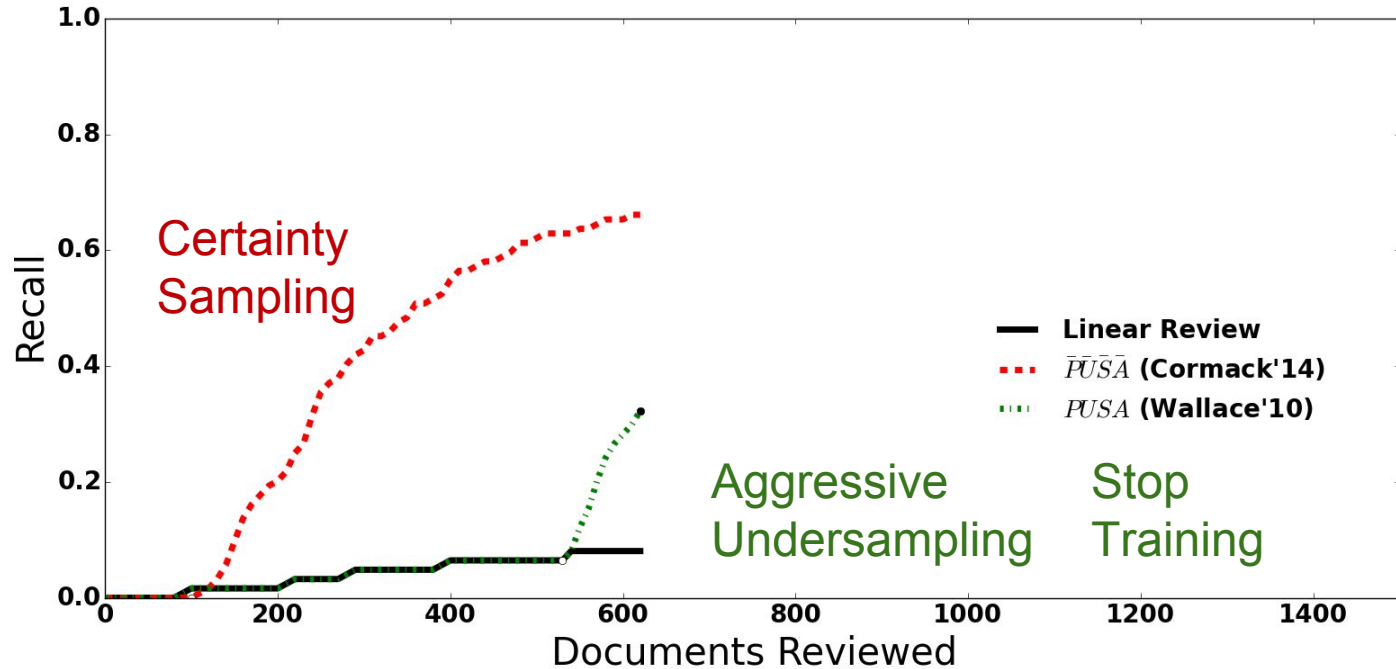
Evaluation 1



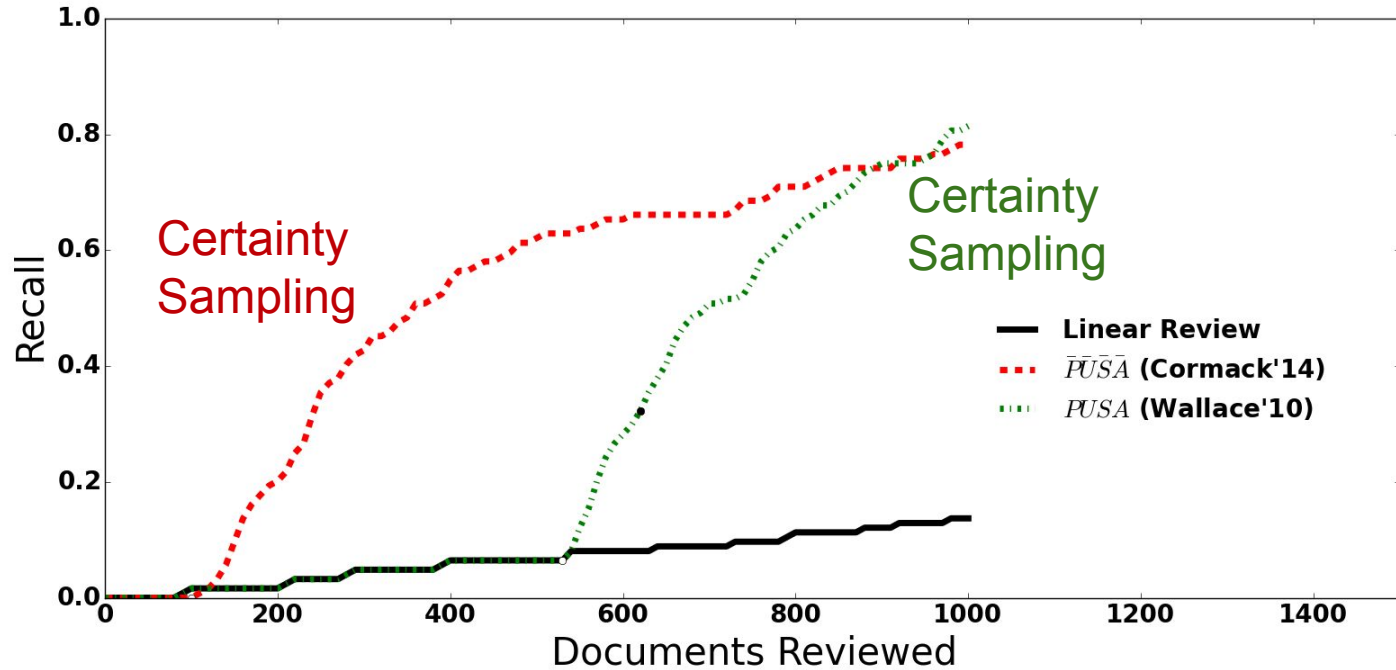
Evaluation 1



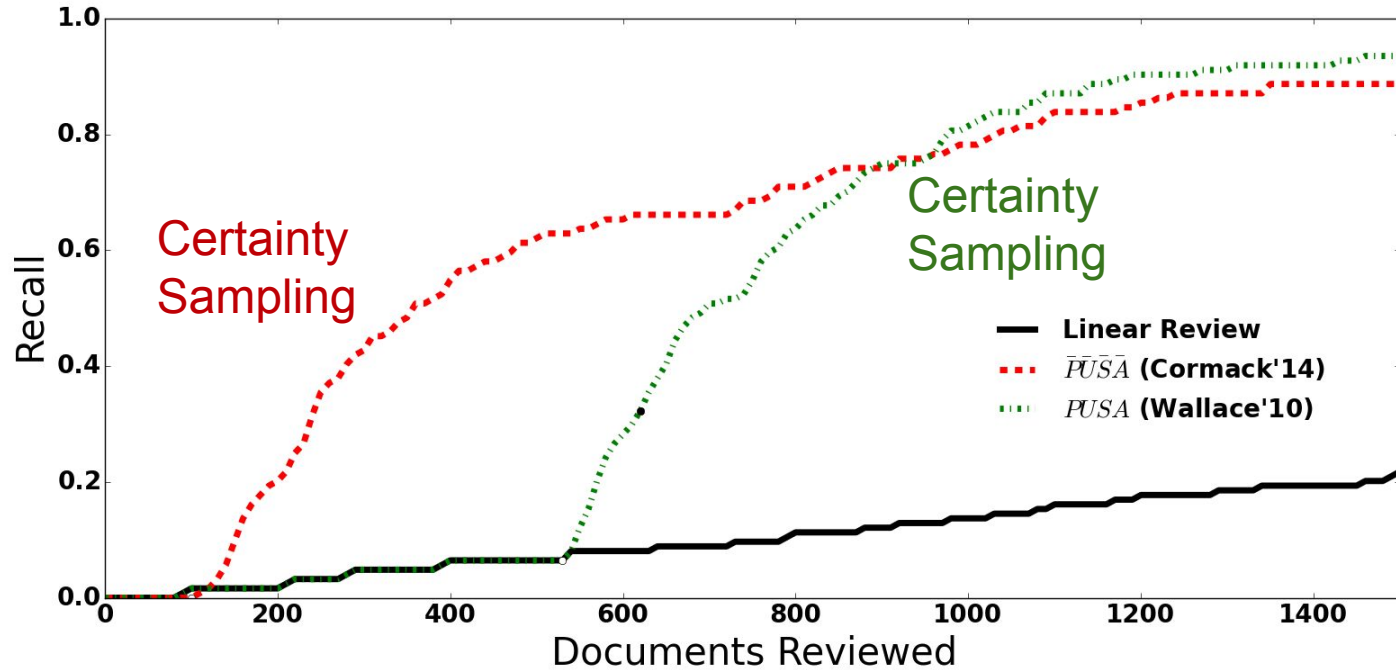
Evaluation 1



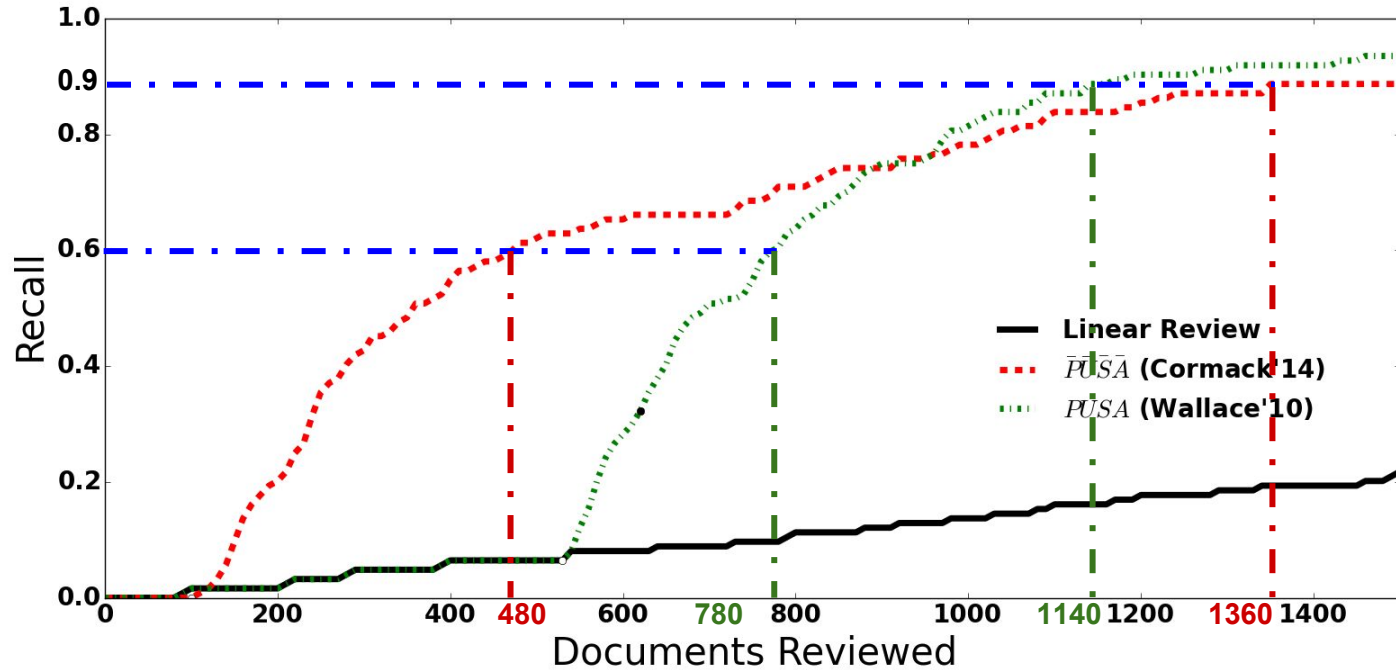
Evaluation 1



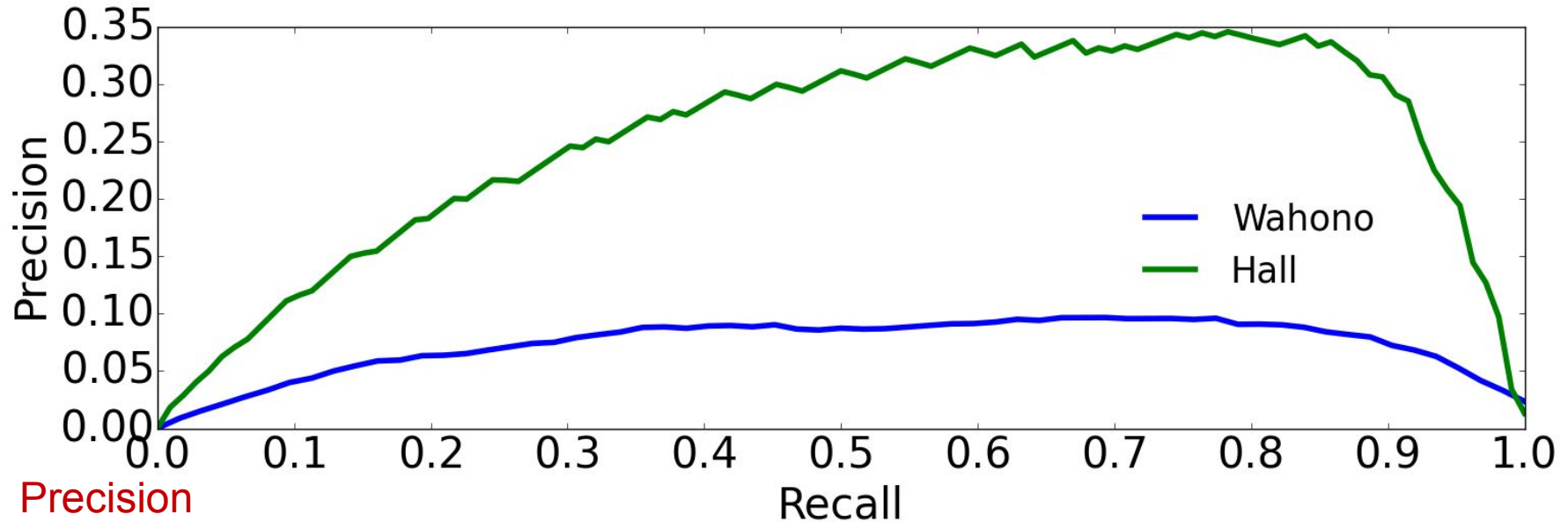
Evaluation 1



Evaluation 1



When to Stop

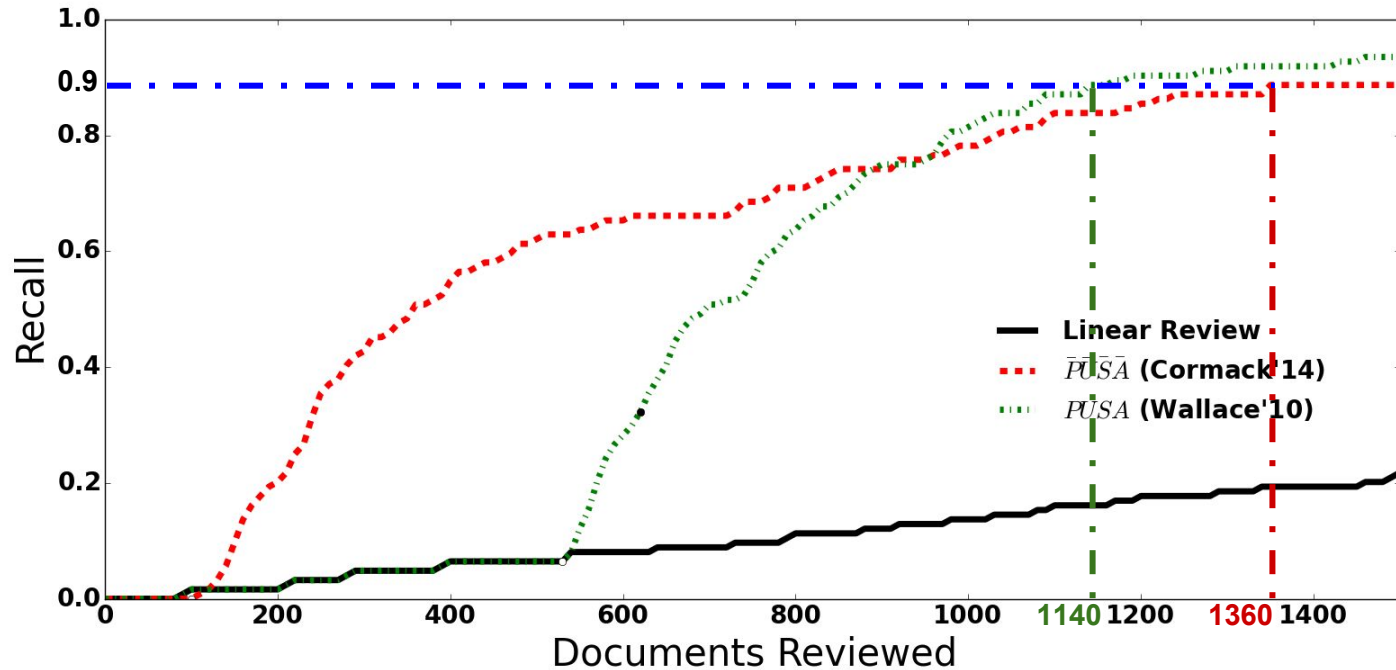


Precision

= relevant found / studies reviewed

Evaluation 2

X90 = studies reviewed to reach 90% recall



Research Questions

RQ1:

Can human-in-the-loop methods improve SE SLR?

RQ2:

Should we adopt state-of-the-art?

RQ3:

How much can we save?

Research Questions

RQ1:

Can human-in-the-loop methods improve SE SLR?

RQ2:

Should we adopt state-of-the-art?

RQ3:

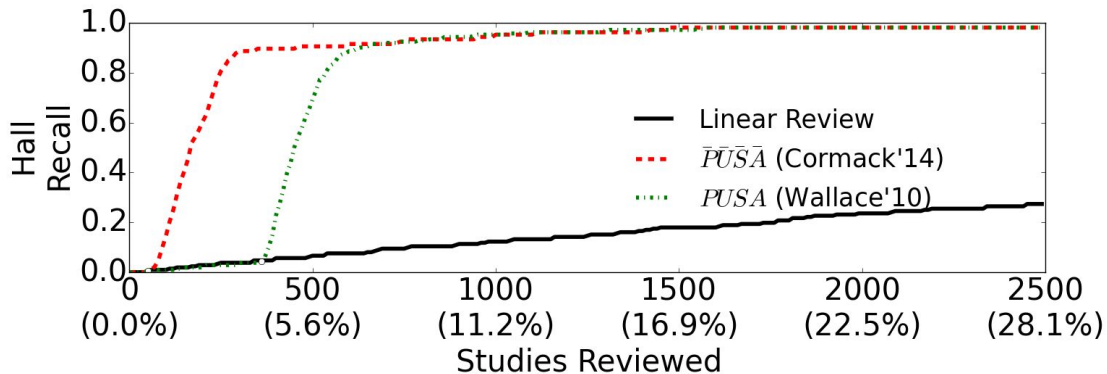
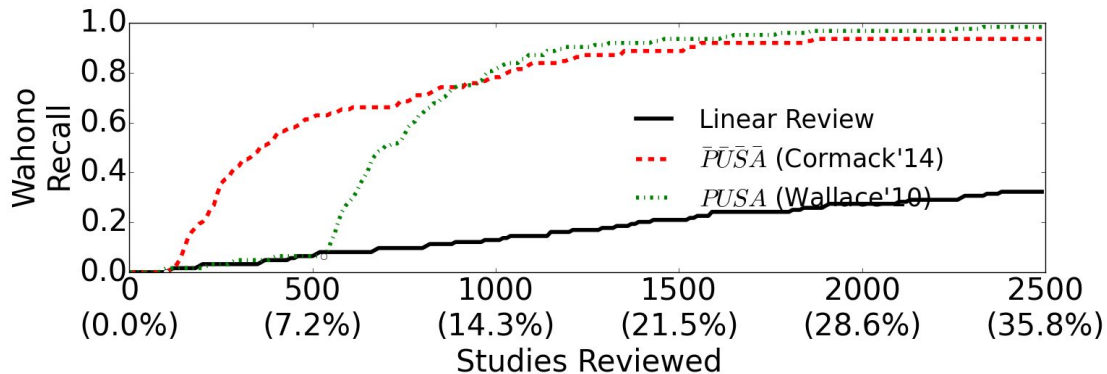
How much can we save?

Outline

- *Background*
- *Method*
- *Experiment*
- ***Result***
 - Human-in-the-loop vs. linear review
 - FASTREAD vs. state-of-the-art
 - Cost Reduction
- *Conclusion and Future works*



Human-in-the-loop vs. Linear Review



Research Questions

RQ1:

Can human-in-the-loop methods improve SE SLR? Yes

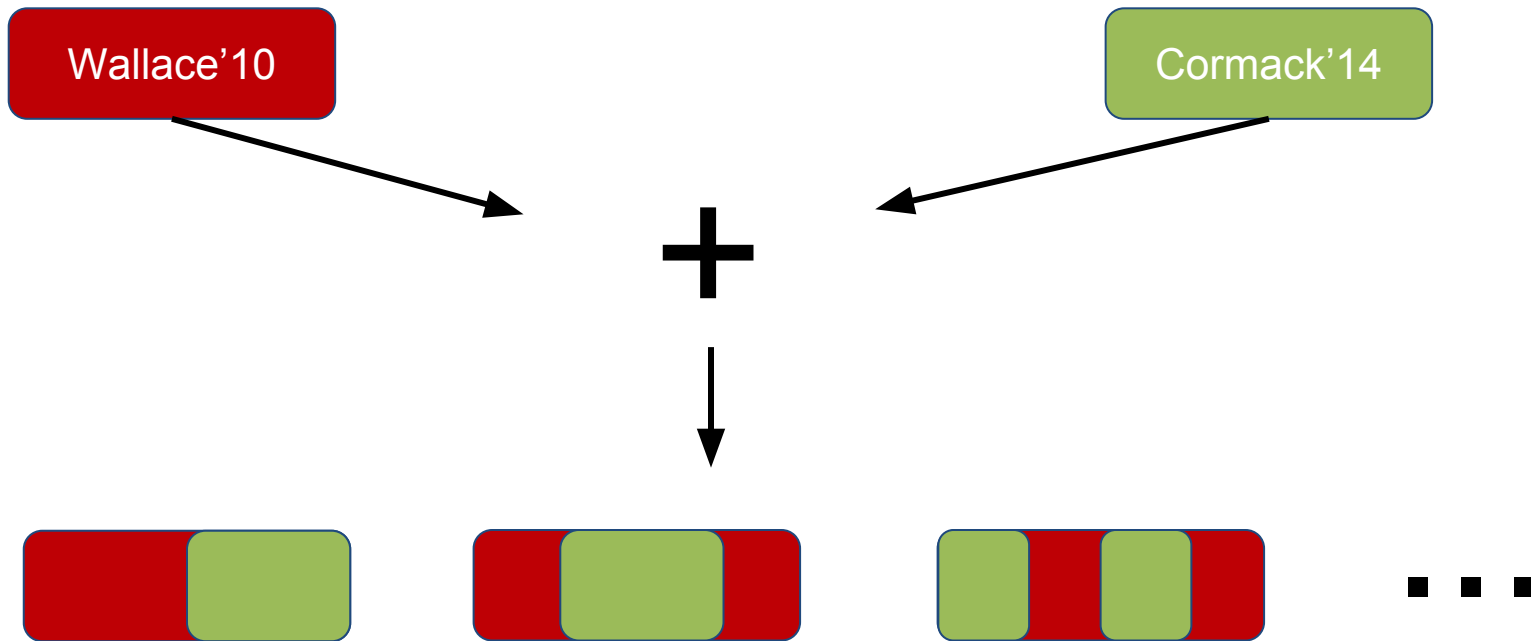
RQ2:

Should we adopt state-of-the-art?

RQ3:

How much can we save?

Refactor



Refactor

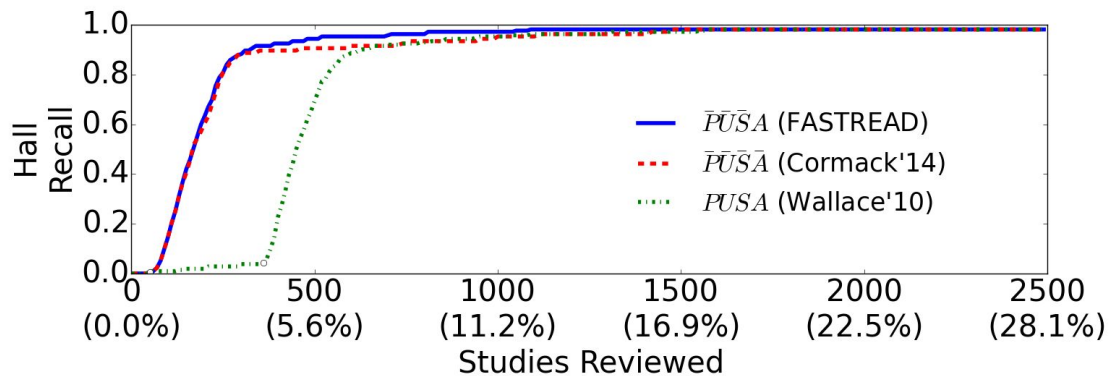
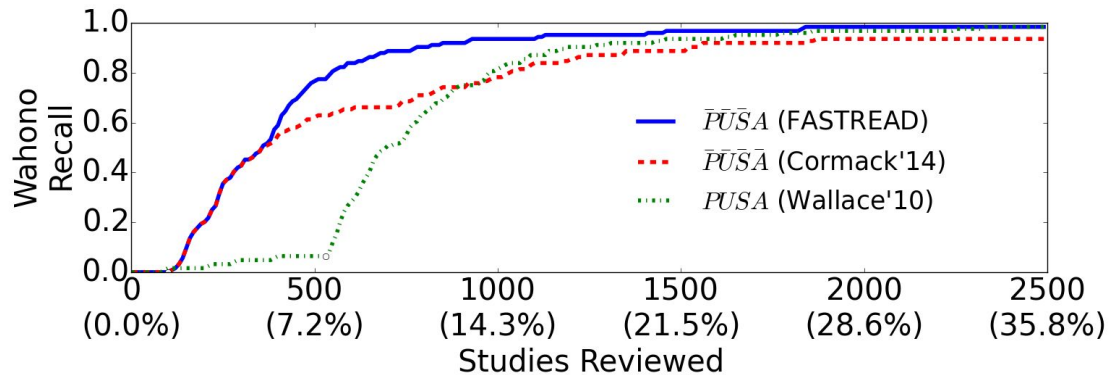
	Start point	Query Strategy	Continuity	Data Balancing
Wallace'10	late (P)	uncertainty sampling (U)	stop training (S)	aggressive undersampling (A)
Cormack'14	early (\bar{P})	certainty sampling (\bar{U})	non-stop (\bar{S})	none (\bar{A})
FASTREAD	early (\bar{P})	certainty sampling (\bar{U})	non-stop (\bar{S})	aggressive undersampling (A)

Outline

- *Background*
- *Method*
- *Experiment*
- ***Result***
 - Human-in-the-loop vs. linear review
 - FASTREAD vs. state-of-the-art
 - Cost Reduction
- *Conclusion and Future works*



FASTREAD vs. State-of-the-art



Research Questions

RQ1:

Can human-in-the-loop methods improve SE SLR? Yes

RQ2:

Should we adopt state-of-the-art? No

RQ3:

How much can we save?

Outline

- *Background*
- *Method*
- *Experiment*
- ***Result***
 - Human-in-the-loop vs. linear review
 - FASTREAD vs. state-of-the-art
 - Cost Reduction
- *Conclusion and Future works*



X90 = studies reviewed
to reach 90% recall

Wahono:

Review 680 (9.7%) to
retrieve 90% relevant

Hall:

Review 310 (3.5%) to
retrieve 90% relevant

X90 from 30 experiments

Rank	Treatment	Median	IQR	Wahono
1	<i>PŪSĀ</i> (FASTREAD)	68C	13C	●
1	<i>PŪSĀ</i>	700	150	●
2	<i>PŪSĀ</i>	840	350	●
2	<i>PŪSĀ</i>	1060	320	●
2	<i>PŪSĀ</i>	1090	330	●
2	<i>PUSĀ</i> (Wallace'10)	114C	28C	●
2	<i>PUSĀ</i>	1340	140	●
2	<i>PŪSĀ</i> (Cormack'14)	135C	13C	●
3	<i>PŪSĀ</i>	1640	370	●
3	<i>PŪSĀ</i>	1640	380	●
4	<i>PŪSĀ</i>	2240	1440	●
4	<i>PUSĀ</i>	2490	770	●
5	Linear Review	624C	31C	●

Rank	Treatment	Median	IQR	Hall
1	<i>PŪSĀ</i>	290	70	●
1	<i>PŪSĀ</i>	300	70	●
1	<i>PŪSĀ</i> (FASTREAD)	31C	3C	●
2	<i>PUSĀ</i>	320	80	●
2	<i>PŪSĀ</i> (Cormack'14)	35C	7C	●
2	<i>PUSĀ</i>	350	100	●
3	<i>PUSĀ</i>	620	260	●
3	<i>PUSĀ</i> (Wallace'10)	63C	27C	●
3	<i>PŪSĀ</i>	660	340	●
3	<i>PŪSĀ</i>	680	240	●
3	<i>PŪSĀ</i>	680	230	●
3	<i>PUSĀ</i>	750	230	●
4	Linear Review	788C	25C	●

Resources

Data:

- [SeaCraft Zenodo](#)

Tool:

- [SeaCraft Zenodo](#)
- [Github](#)

The screenshot shows a web application interface for document classification and analysis. The interface is divided into several sections:

- File Selection (A):** A 'Choose File' button is highlighted with a red arrow. Below it, 'Next' and 'Auto Review' buttons are visible.
- Document List (B):** A list of documents is displayed. The sixth item, '6. Estimating software fault-proneness for tuning testing activities (0.5667289538574896)', is highlighted in yellow and has a red arrow pointing to it.
- Classification Control (C):** A 'Random' button is highlighted with a red arrow. Below it, 'Certain' and 'Uncertain' radio buttons are visible.
- Submit (D):** A 'Submit' button is highlighted with a red arrow.
- Graph (E, F, G, H, I, J):** A line graph titled 'Documents Coded: 92/280 (8911)' shows 'Relevant Found' on the y-axis (0 to 100) and 'Documents Reviewed' on the x-axis (0 to 300). The graph shows a blue line that starts at (0,0) and rises to approximately (300, 90). Red arrows point to specific points on the graph: H at (0,0), I at (100,0), and J at (150,0).

Outline

- *Motivation and Background*
- *Method*
- *Experiment*
- *Result*
- ***Conclusion and Future works***



Conclusion

- *Apply human-in-the-loop method to facilitate SLRs*
- *FASTREAD, a better method than state-of-the-art*
- *Save 90% cost to retrieve 90% relevant studies*
- *A tool to implement FASTREAD*

Future Roadmap

Assumptions:

- no external domain knowledge
- binary classification
- one reviewer, no error

Future Roadmap

Assumptions:

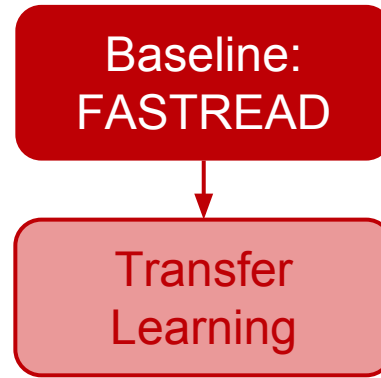
Baseline:
FASTREAD

- no external domain knowledge
- binary classification
- one reviewer, no error

Future Roadmap

Assumptions:

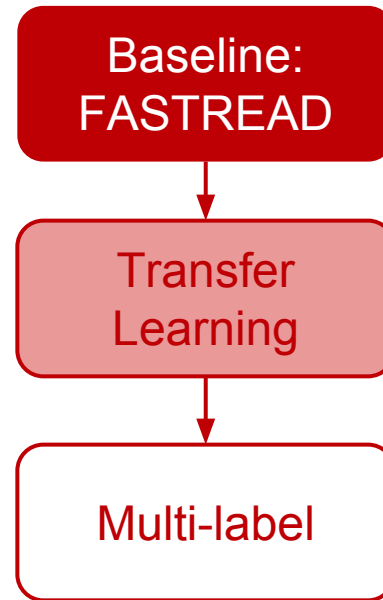
- no external domain knowledge
- binary classification
- one reviewer, no error



Future Roadmap

Assumptions:

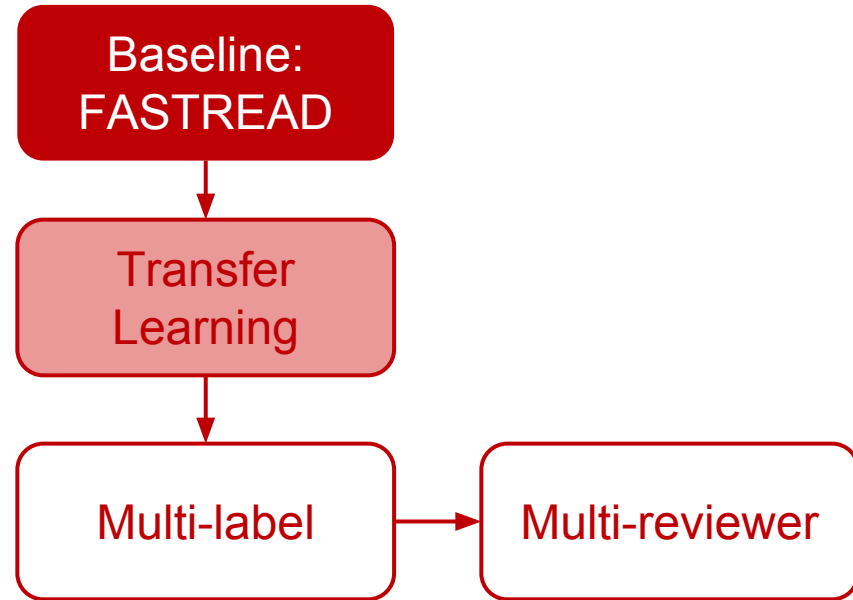
- no external domain knowledge
- binary classification
- one reviewer, no error



Future Roadmap

Assumptions:

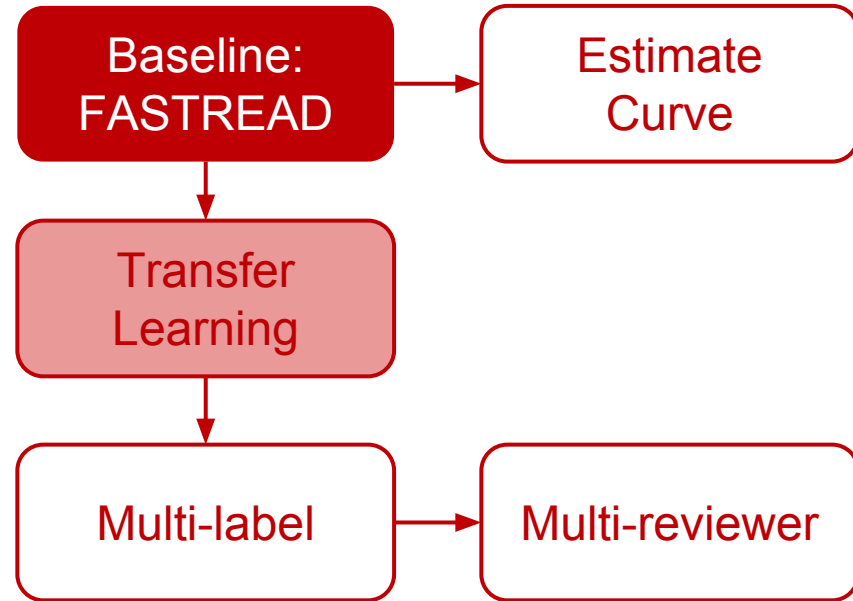
- no external domain knowledge
- binary classification
- one reviewer, no error



Future Roadmap

Assumptions:

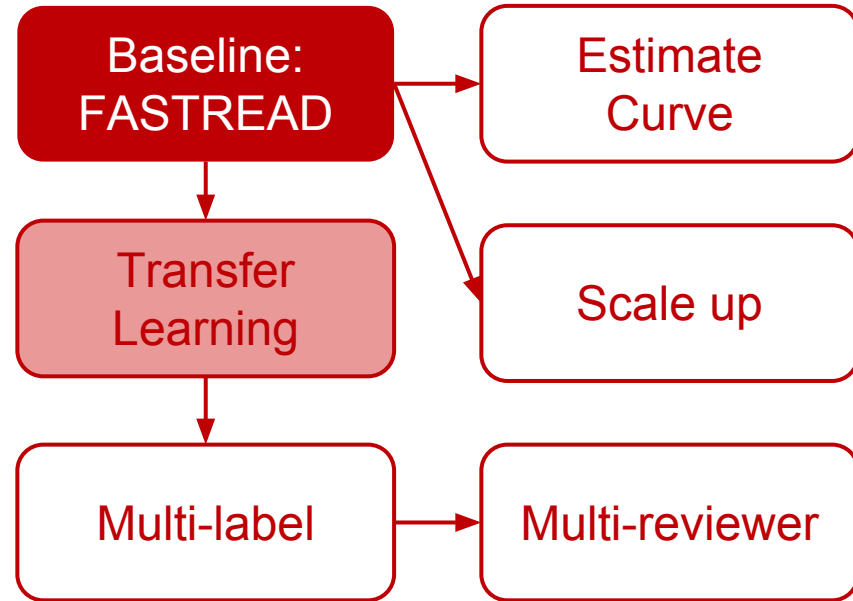
- no external domain knowledge
- binary classification
- one reviewer, no error



Future Roadmap

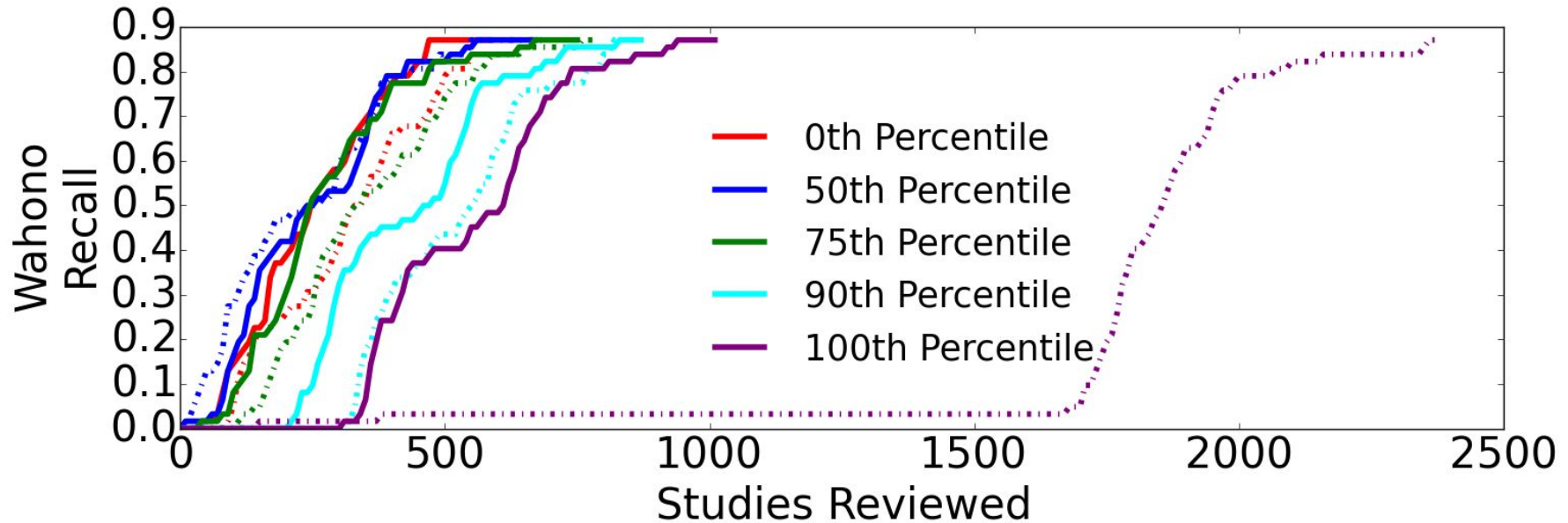
Assumptions:

- no external domain knowledge
- binary classification
- one reviewer, no error

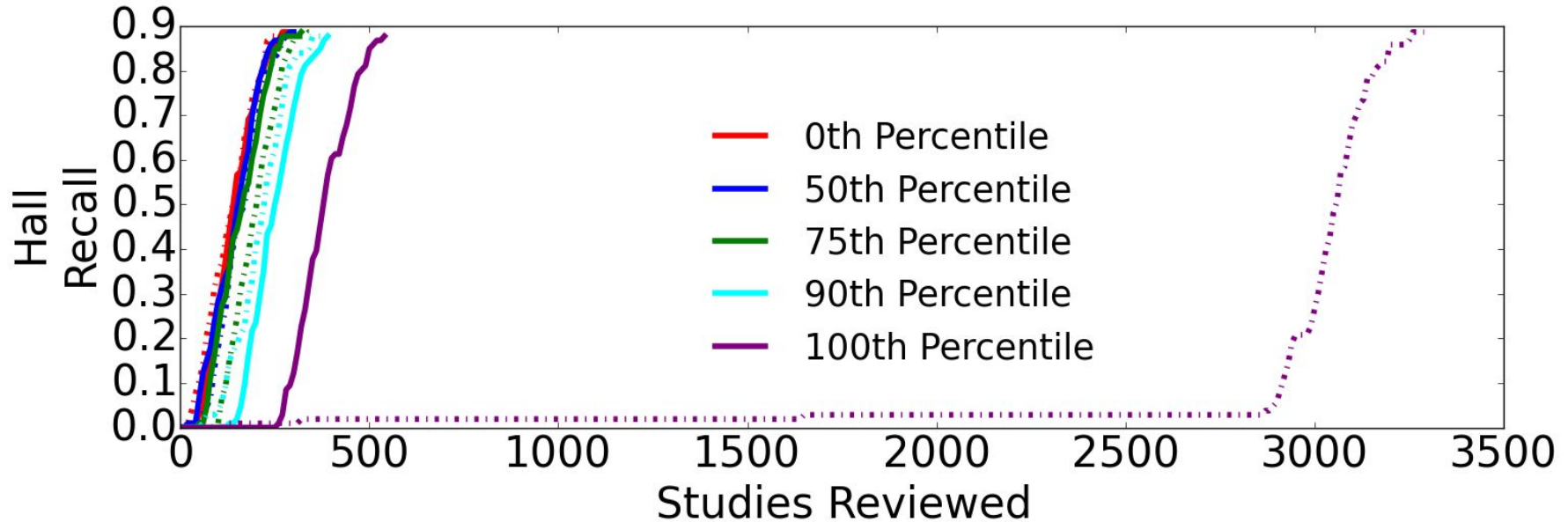




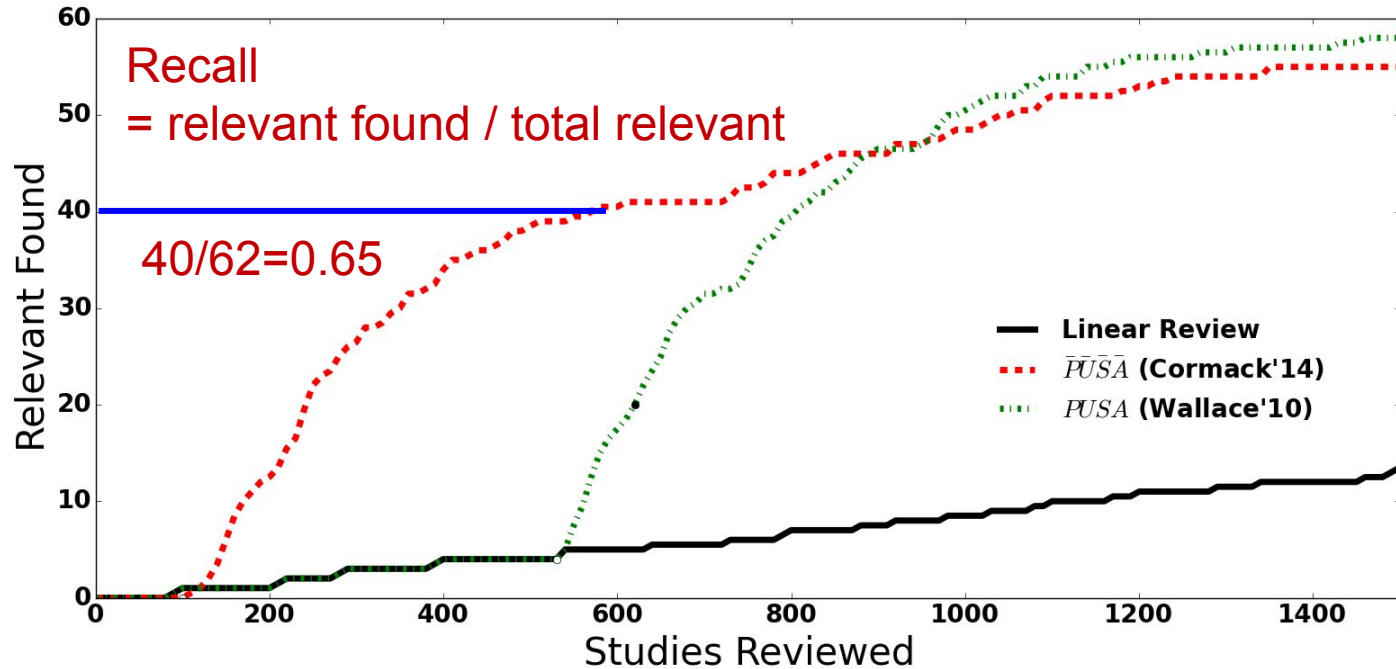
Back-up Slides



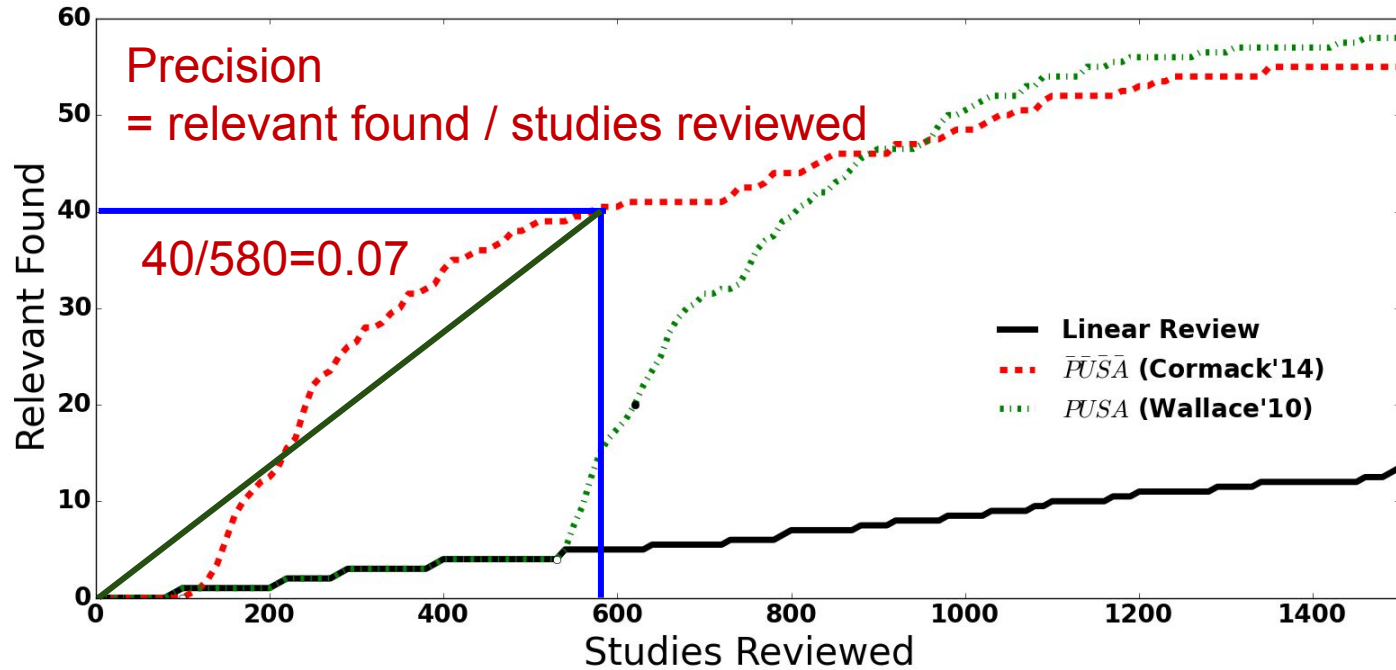
Back-up Slides



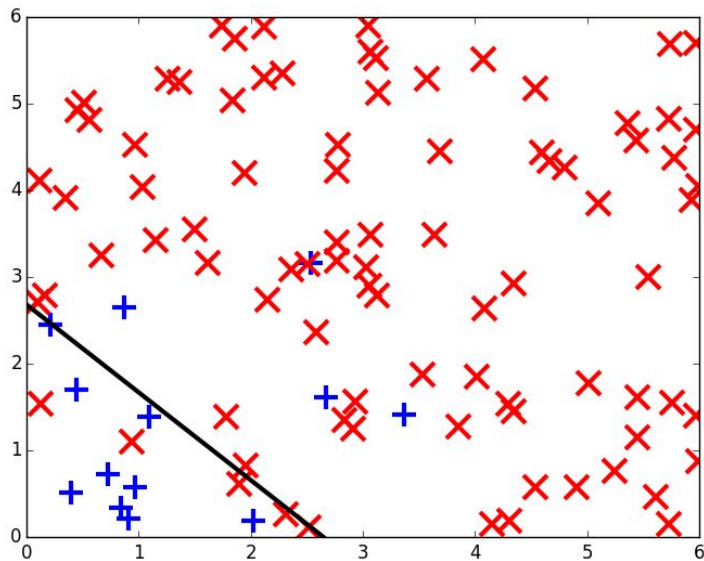
Evaluation



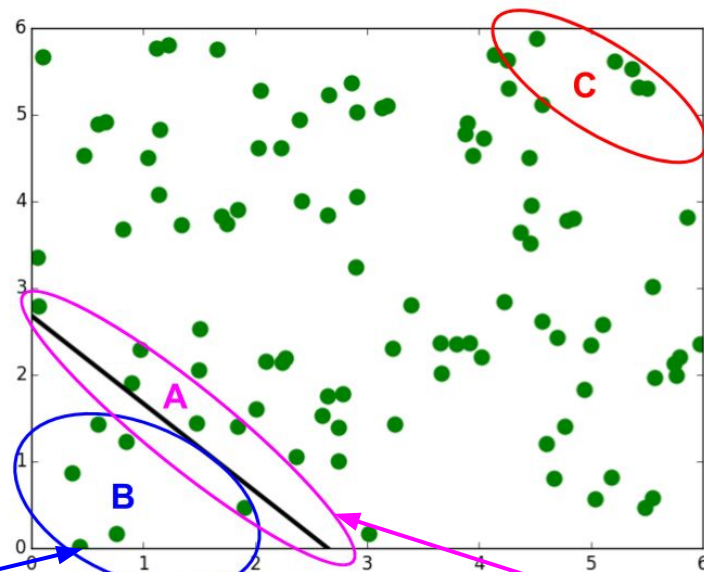
Evaluation



Query Strategy

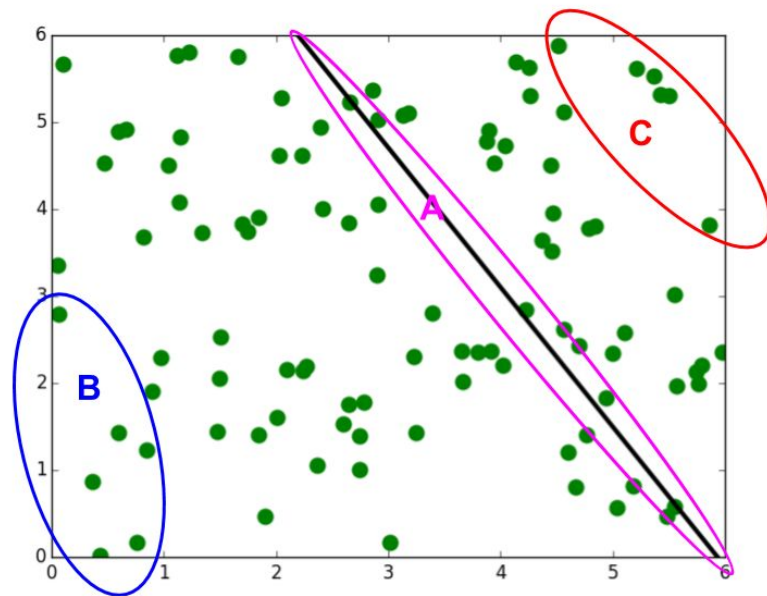
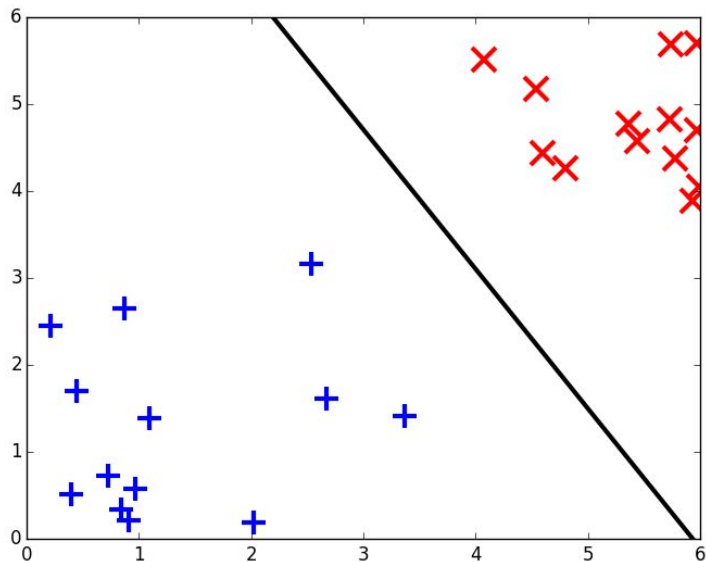


Certainty
Sampling



Uncertainty
Sampling

Data Balancing



Data Sets

	Wahono		Hall	
	Stated	Retrieved	Stated	Retrieved
Initial List	2117	7002	2073	8911
Final List	72	62	136	106

Roadmap

Systematic literature reviews:

- Useful. Important.
 - In medicine, law, SE and elsewhere
- Expensive.

Cost Reduction

- Primary study selection= part of each reviews. Hard.
- Tools to make it easier?
 - Problems with state of the art

Human-in-the-loop Incremental Learning

- Used in other domains.
- Q: Useful in SE?
- A1: Yes.
- A2: Can be used to tailor a even better SE method.





Documents Ranked

Documents Coded: 92/280 (8911)

Relevant | Irrelevant | Undetermined | None

True Label: yes

Estimating software fault-proneness for tuning testing activities

The article investigates whether a correlation exists between the fault-proneness of software and the measurable attributes of the code (i.e. the static metrics) used for testing (i.e. the dynamic metrics). The article also studies how to use such data for tuning the testing process. The goal is not to find a general solution to the problem (a solution may not even exist), but to investigate the scope of specific solutions, i.e., to what extent homogeneity of the development process, organization, environment and application domain allows a data computed on past projects to be projected onto new projects. A suitable variety of case studies is selected to investigate a methodology applicable to classes of heterogeneous products, rather than investigating if a specific solution exists for five cases.

