

Estimating Summary Quality with Pairwise Preferences



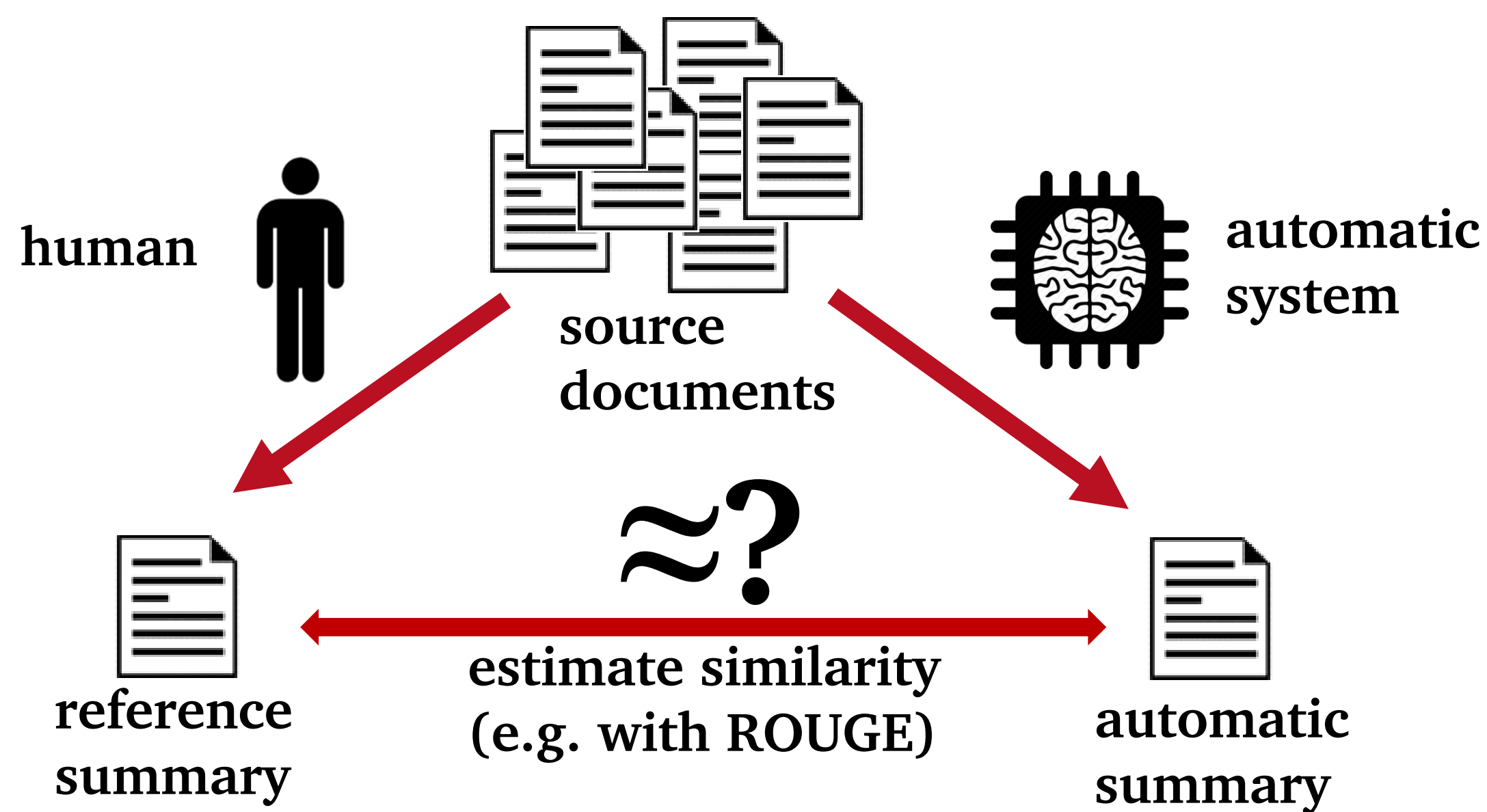
TECHNISCHE
UNIVERSITÄT
DARMSTADT



AIPHES

Markus Zopf | Research Training Group AIPHES, TU Darmstadt

Traditional Automatic Evaluation



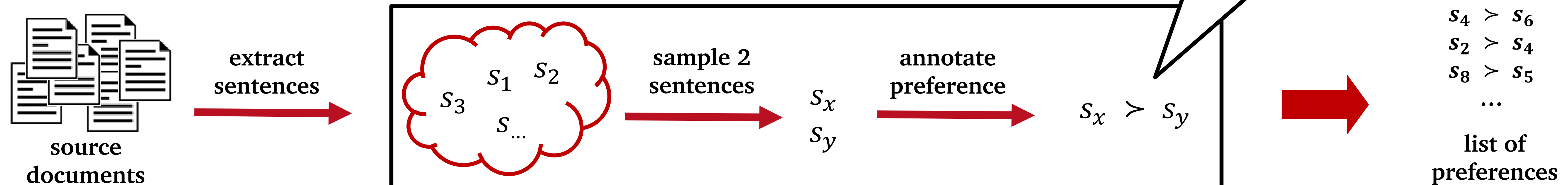
Problems:

1. creating reference summaries is time-consuming and complex → expensive ⚡
2. computing text similarities is a complex, unsolved problem → not reliable ⚡

Summary

- 📍 traditional evaluation methods use references summaries to estimate the quality of automatically generated summaries
- ⚡ creating references summaries is complex and expensive
- 💡 we propose to use simple pairwise preference labels of individual sentences for evaluation
- 💡 our evaluation is cheaper and more accurate than the SOTA

Estimating Summary Quality with Pairwise Preferences



generated preferences are interpreted as games in which better sentences are more likely to win a game

$$s_4 > s_6$$

$$s_2 > s_4$$

$$s_8 > s_5$$

$$\dots$$

compute Bradley-Terry scores

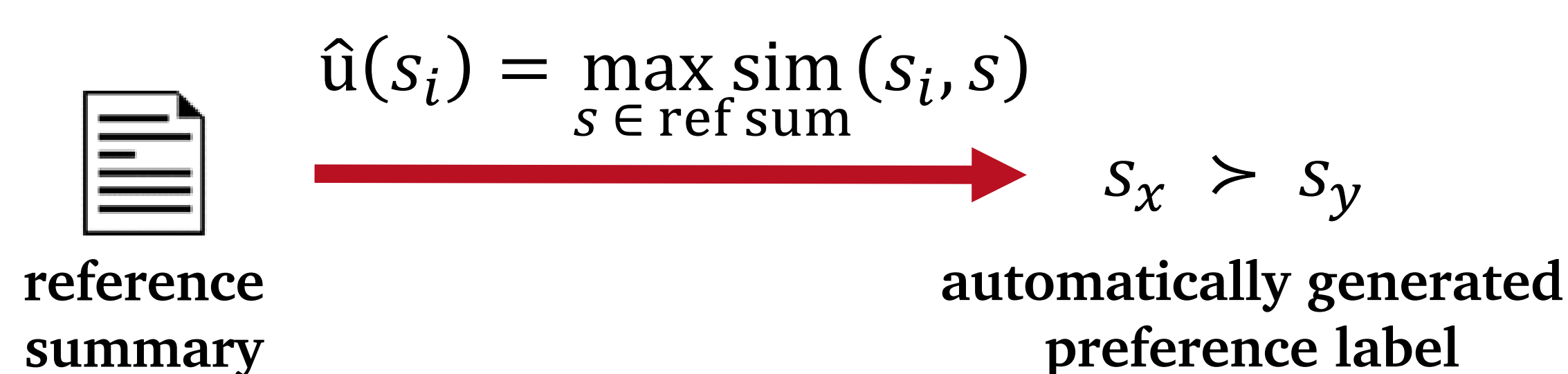
$$p(s_x > s_y) = \frac{u(s_x)}{u(s_x) + u(s_y)}$$

$$v(\mathbf{S}) = \sum_i^{|\mathbf{S}|} w_i \cdot u(s_i)$$

\mathbf{S} = summary to evaluate
 s_i = i -th sentences in \mathbf{S}

Automatically Generating Preferences

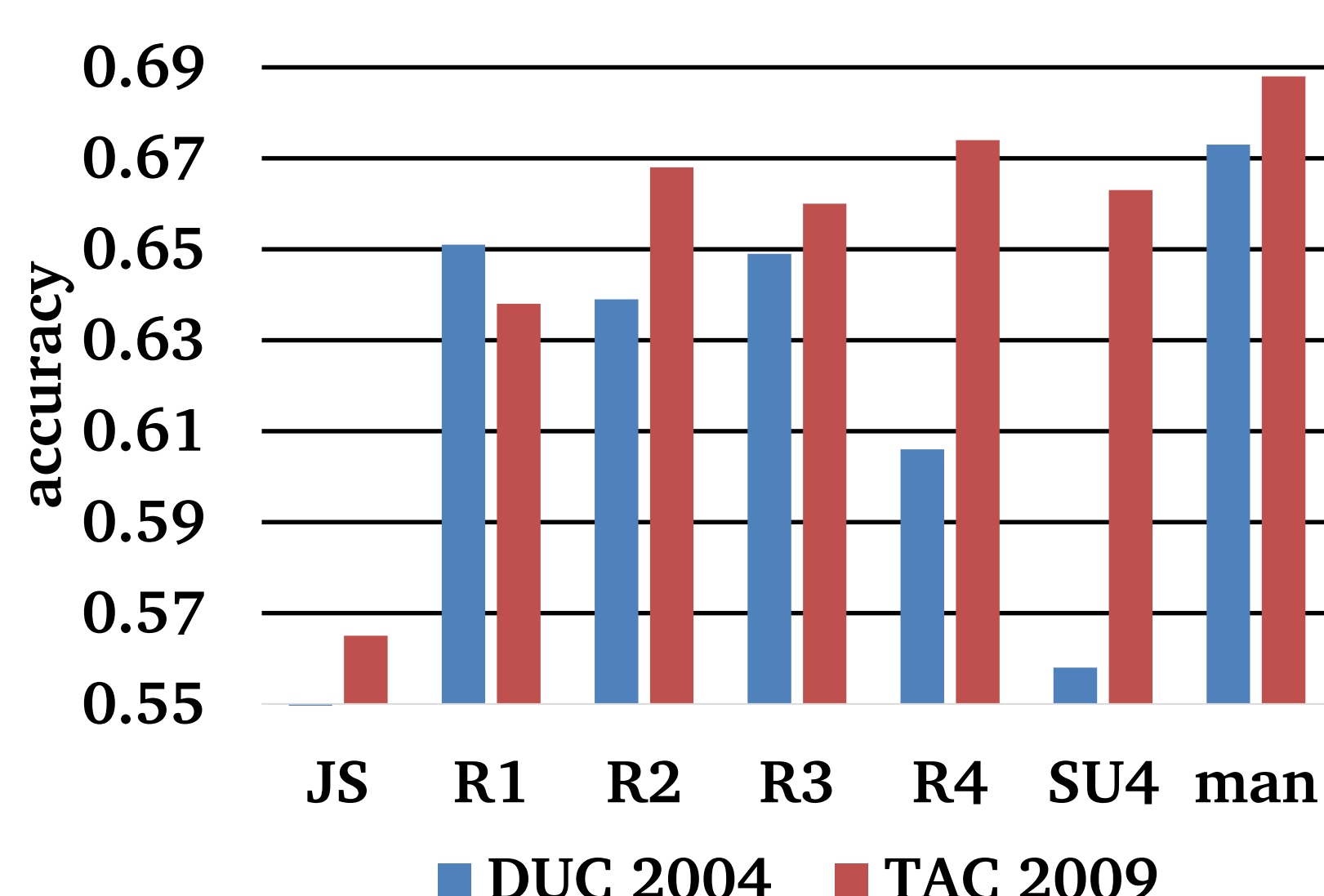
reference summaries are already available for standard datasets → reuse reference summaries to automatically create preference annotations for free



Evaluating Evaluation Methods

<p>Prior work: Pearson correlation</p> <p>interpretation difficult</p> <p>requires linear correlation</p> <p>normal distr., interval scaling assumed</p> <p>sensible to outliers</p>	<p>Our work: Pairwise accuracy</p> <p>clear interpretation</p> <p>only necessary requirement</p> <p>no additional assumptions</p> <p>robust</p>
--	---

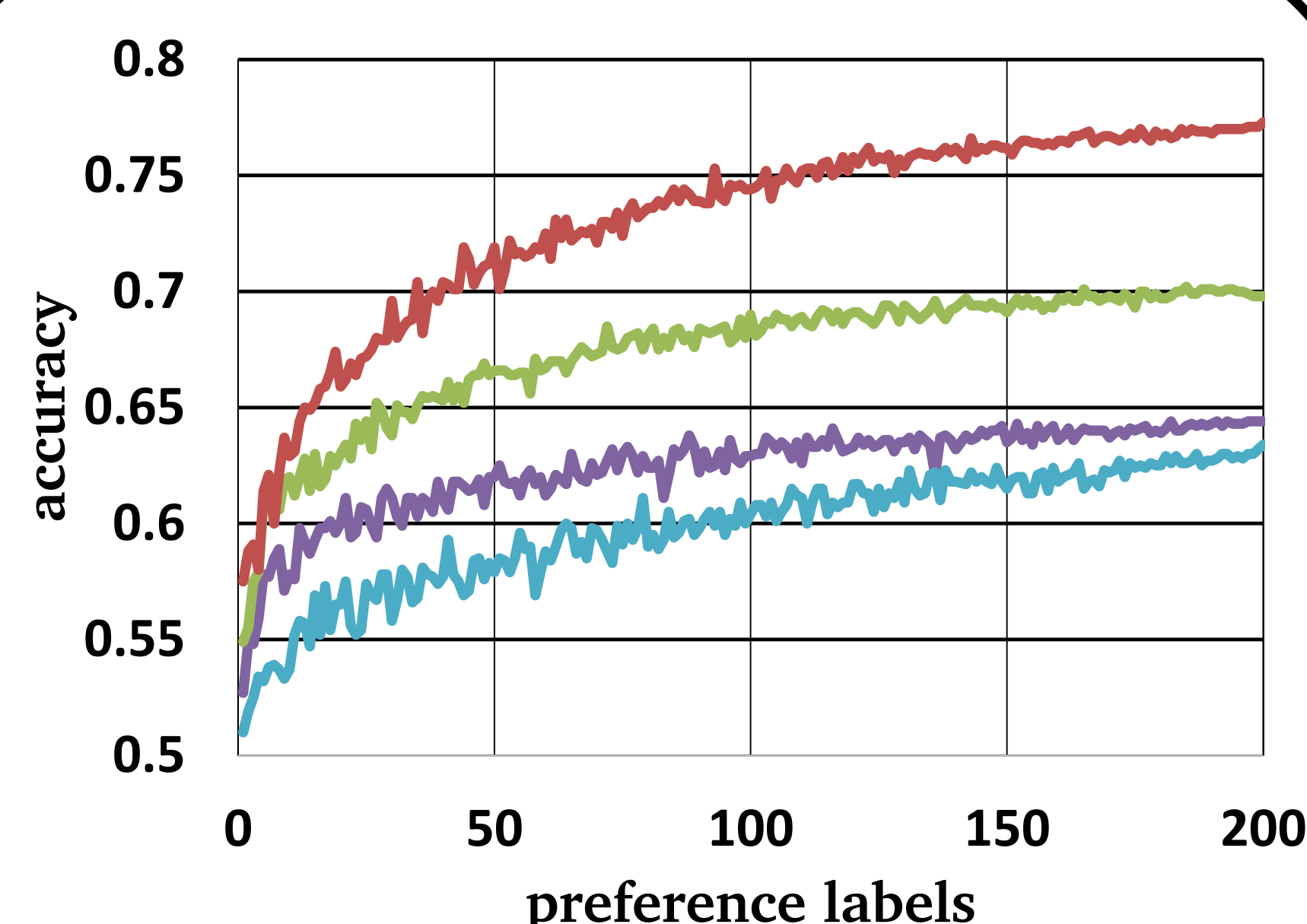
Using 200 Human Preferences



Observations:

- annotation effort: ~54 minutes per topic
- preference-based evaluation with human annotations outperforms JS and ROUGE

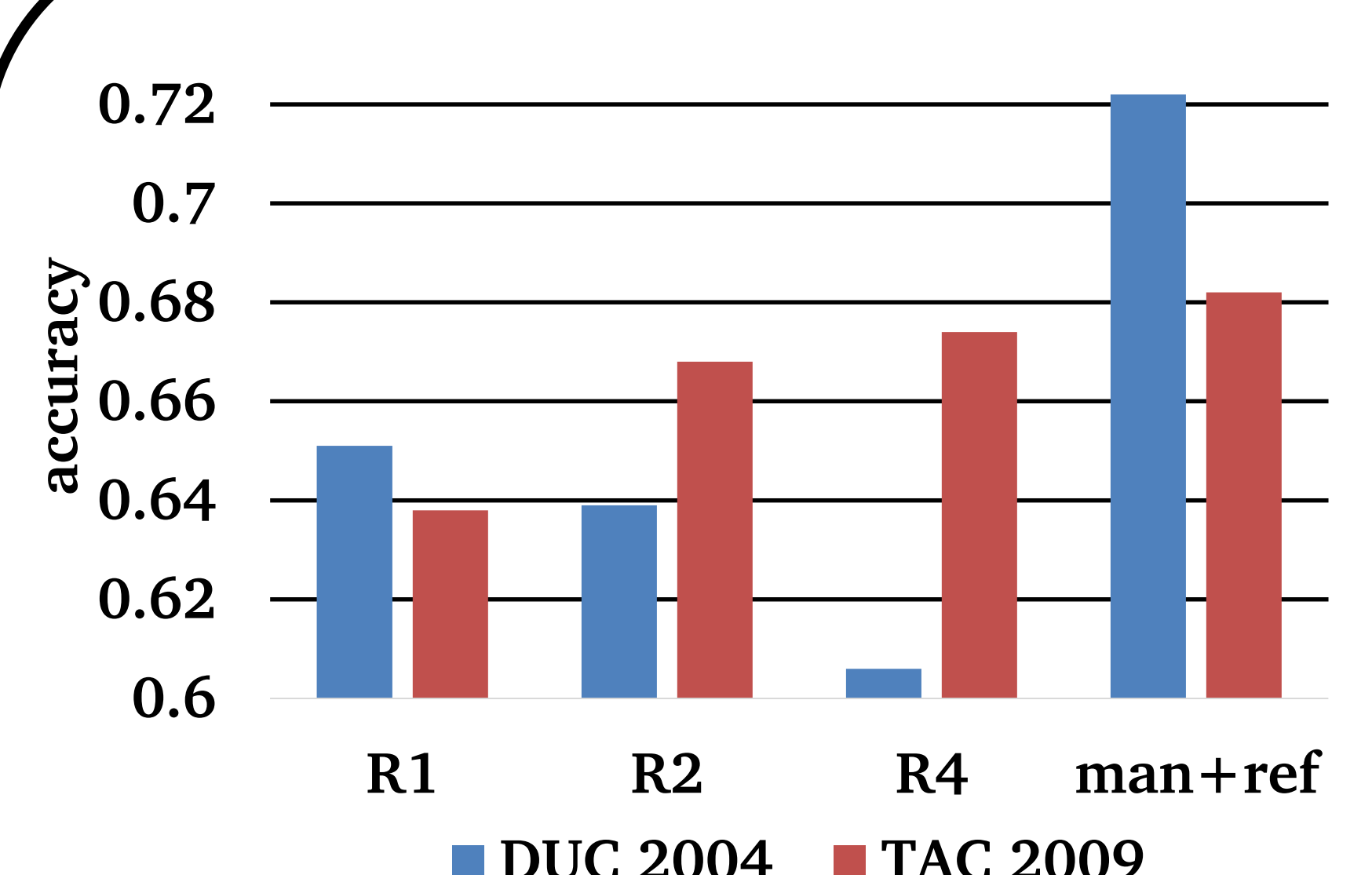
Training Convergence



Observations:

- different topics → different convergence
- not yet converged after 200 labels

Human+Automatic Preferences



Observation:

- using additional automatically generated labels can further improve accuracy