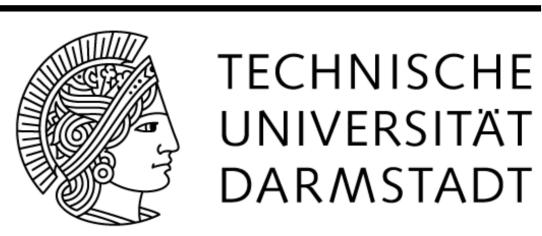
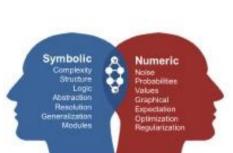
## STRICTA: Structured Reasoning In Critical Text Assessment for Peer Review and Beyond







Nils Dycke, Matej Zečević, Ilia Kuznetsov, Beatrix Suess, Kristian Kersting, Iryna Gurevych

UKP Lab, AIML Lab, Synthetic Biology Lab at Technical University of Darmstadt



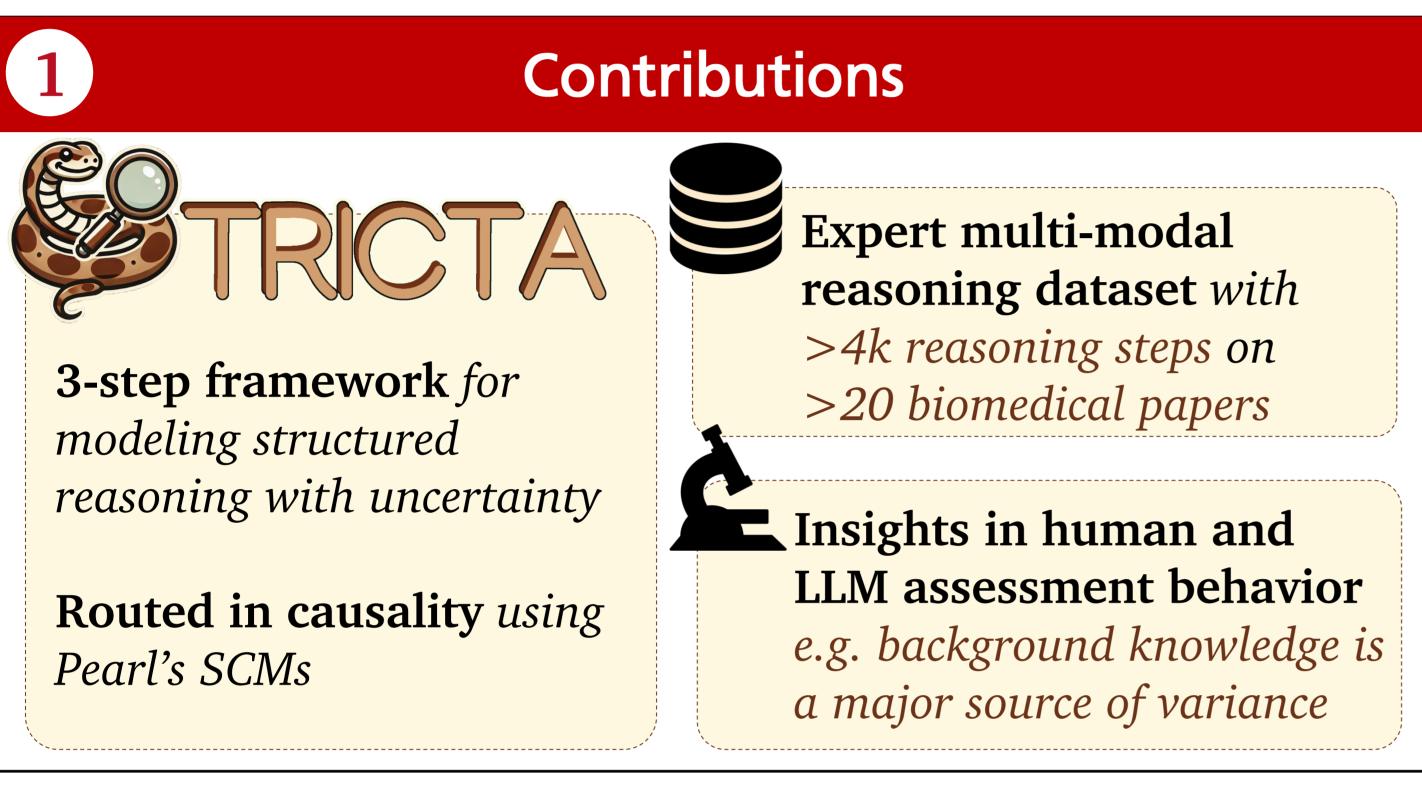


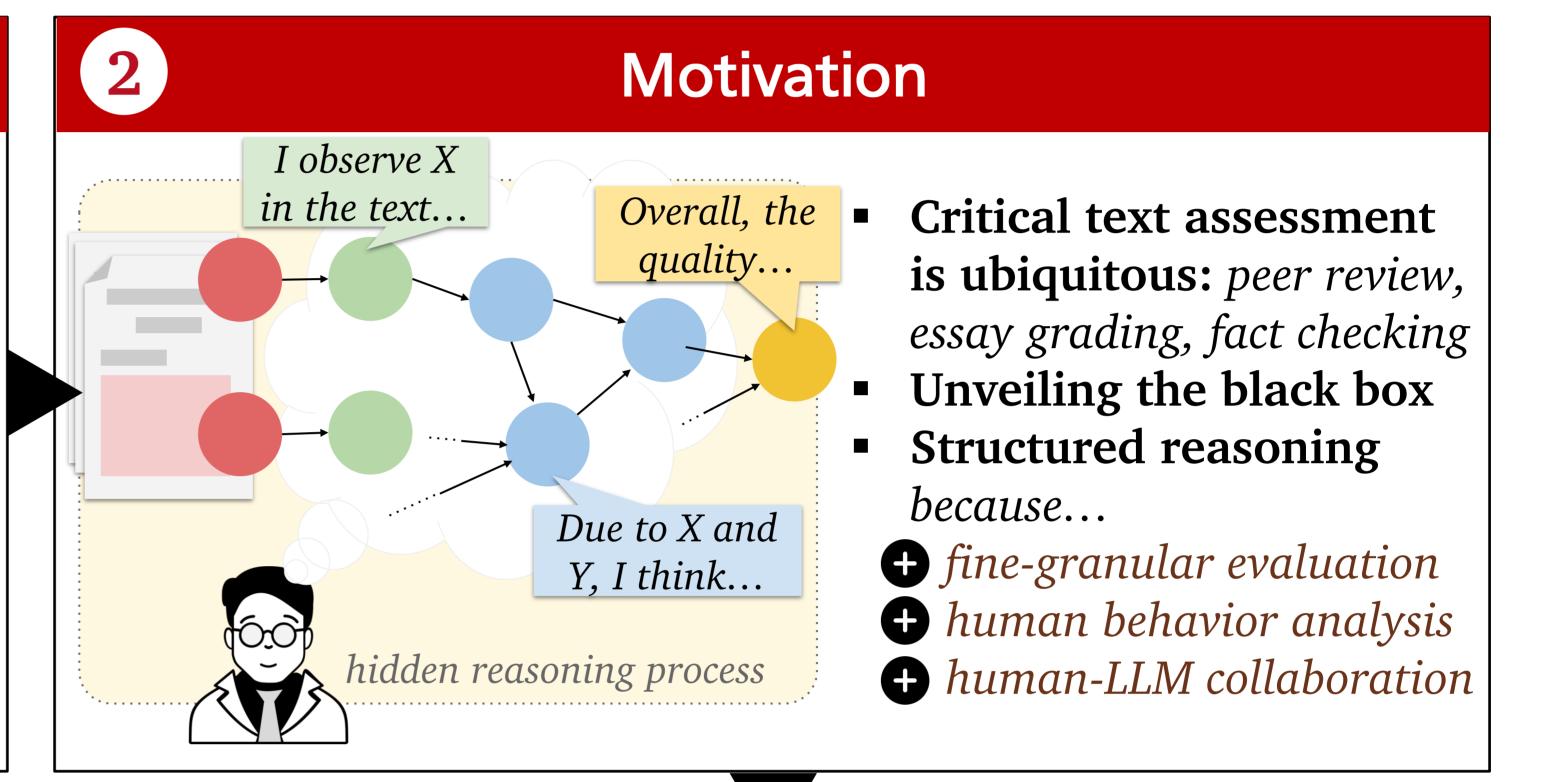


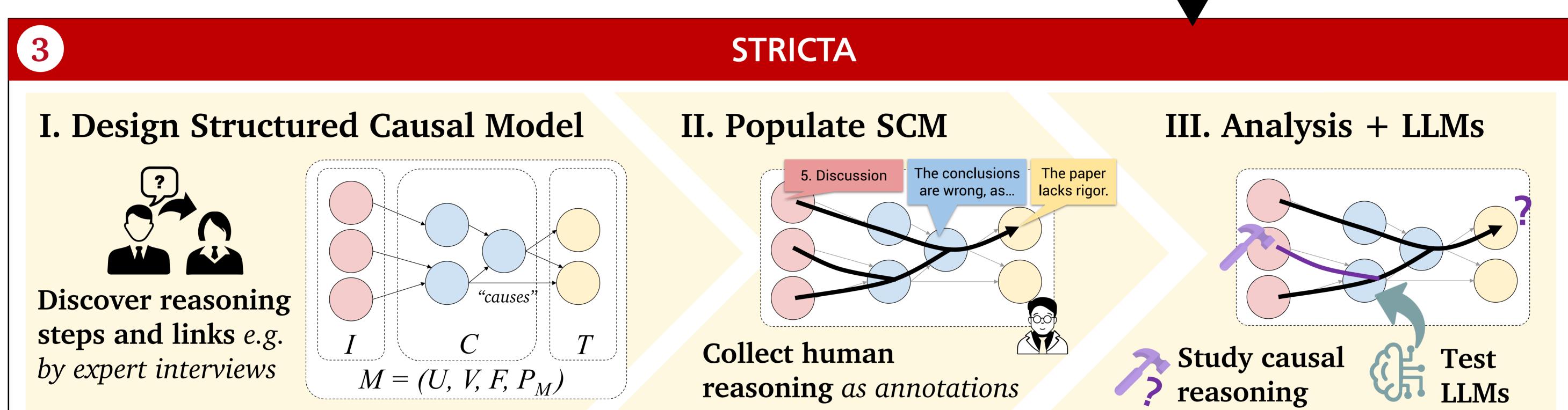


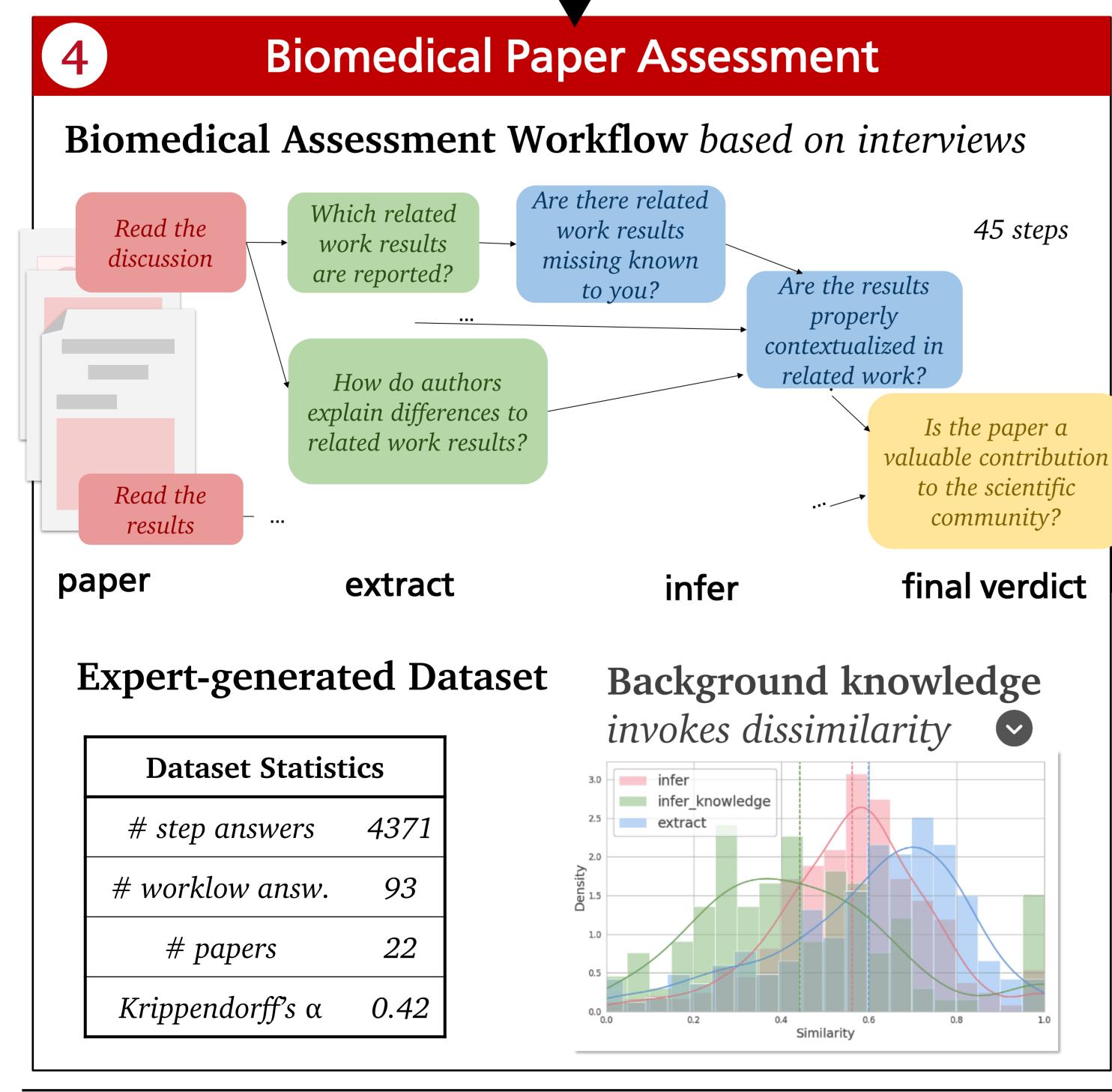












## Causal Analysis + LLM Experiments

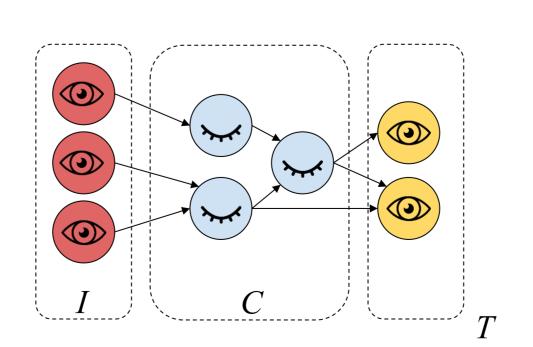
## Which factors shape the assessment most?

Steps' Average Causal Effect on Verdict				
Consistency of conclusions to RQs	0.37			
Relevance of conclusions to science	0.20			
Clarity of the paper	0.20			

Humans show a positive bias towards well-written papers.

... and a lot more analysis including counterfactuals etc.

## Can LLMs explain human's final verdicts?



LLMs perform best on extraction and weight factors differently

	BERT-F1↑	$\mathbf{SummaC} \uparrow$	<b>TRUE</b> ↑	<b>F1</b> ↑
human*†	$\boldsymbol{0.799}_{\pm.06}$	$-0.151_{\pm .30}$	$0.151_{\pm .27}$	0.801
Llama3 Prg <sup>†</sup>	$0.752_{\pm .10}$	$-0.274_{\pm .36}$	$0.098_{\pm .30}$	0.170
Mixtral Prg <sup>†</sup>	$0.761_{\pm .09}$	$ extbf{-0.149}_{\pm .26}$	$0.120_{\pm .32}$	0.559
GPT3.5t Prg <sup>†</sup>	$0.759_{\pm .10}$	$-0.178_{\pm .35}$	$\boldsymbol{0.163}_{\pm .37}$	0.531
GPT4o Prg <sup>†</sup>	$0.780_{\pm .07}$	$-0.186_{\pm .30}$	$0.139_{\pm .35}$	0.720
majority <sup>io</sup>				0.854
human $^{*io}$	$0.799_{\pm .06}$	$-0.158_{\pm .29}$	$0.150_{\pm .27}$	0.801
Llama3 <sup>io</sup>	$0.786_{\pm .07}$	$-0.141_{\pm .30}$	$0.145_{\pm .35}$	0.657
$Mixtral^{io}$	$0.794_{\pm .07}$	-0.077 $_{\pm.27}$	$0.161_{\pm .37}$	0.822
GPT3.5t <sup>io</sup>	$\boldsymbol{0.805}_{\pm.07}$	$-0.125_{\pm .34}$	$\textbf{0.214}_{\pm.41}$	0.789
GPT4o <sup>io</sup>	$0.795_{\pm .07}$	$-0.094_{\pm .28}$	$0.194_{\pm .40}$	0.876
GPT4o§	$0.776_{\pm .07}$	$-0.154_{\pm0.28}$	$0.188_{\pm .39}$	0.828









