

Task	Motivation	BioRelFact Dataset																																																														
<p>Joint Detection of Relation & Epistemic Commitment</p> <p>Given a statement in a sentence with two highlighted entities, predict:</p> <p>Relation: Which relation is discussed?</p> <p>Epistemic Commitment: How committed is the statement author to the factual state of the relation?</p> <table border="1"> <thead> <tr> <th colspan="2">Relation: Therapeutic_Use</th> </tr> <tr> <th>Factuality</th> <th>Example</th> </tr> </thead> <tbody> <tr> <td>Factual</td> <td>Drug treats Disease</td> </tr> <tr> <td>Possible</td> <td>Drug may help with Disease</td> </tr> <tr> <td>Doubtful</td> <td>Drug is unlikely to benefit Disease</td> </tr> <tr> <td>Negated</td> <td>Drug does not treat Disease</td> </tr> <tr> <td>Uncommitted</td> <td>Drug is being studied for Disease</td> </tr> </tbody> </table>	Relation: Therapeutic_Use		Factuality	Example	Factual	Drug treats Disease	Possible	Drug may help with Disease	Doubtful	Drug is unlikely to benefit Disease	Negated	Drug does not treat Disease	Uncommitted	Drug is being studied for Disease	<p>Factuality-aware extraction of biomedical relations is crucial for computational drug development</p> <ul style="list-style-type: none"> Standard relation extraction alone has limitations <ul style="list-style-type: none"> → it may either propagate misleading information (if speculative statements are treated as factual) → or lose important nuances (if uncertain or negated relations are discarded). <p>Research Gap & Contributions</p> <ul style="list-style-type: none"> Most RE datasets ignore factuality, e.g., collapsing speculative and factual relations or treating negated relations as unmentioned <ul style="list-style-type: none"> → we publish dataset BioRelFact labeled with relations and epistemic commitment LLMs for factuality-aware biomedical RE are largely unexplored. <ul style="list-style-type: none"> → we test 8 LLMs (commercial vs. open weight, general purpose vs. biomedical) 	<p>Creation of dataset for relation and epistemic commitment detection</p> <ul style="list-style-type: none"> 1,767 sentences from PubMed abstracts Annotated by biomedical experts <table border="1"> <thead> <tr> <th>Relation</th> <th>Count</th> </tr> </thead> <tbody> <tr> <td>no_relation</td> <td>738</td> </tr> <tr> <td>drug : disease</td> <td></td> </tr> <tr> <td>Causal_Effect</td> <td>63</td> </tr> <tr> <td>Therapeutic_Use</td> <td>246</td> </tr> <tr> <td>drug : gene</td> <td></td> </tr> <tr> <td>Agonist</td> <td>22</td> </tr> <tr> <td>Antagonist</td> <td>89</td> </tr> <tr> <td>Modulates</td> <td>42</td> </tr> <tr> <td>gene : disease</td> <td></td> </tr> <tr> <td>Biomarker</td> <td>52</td> </tr> <tr> <td>Causal_Effect</td> <td>47</td> </tr> <tr> <td>Modulates</td> <td>168</td> </tr> <tr> <td>geneVariant : disease</td> <td></td> </tr> <tr> <td>Association</td> <td>300</td> </tr> <tr> <td>Total</td> <td>1767</td> </tr> </tbody> </table> <table border="1"> <thead> <tr> <th>Factuality level</th> <th>Count</th> </tr> </thead> <tbody> <tr> <td>fact</td> <td>832</td> </tr> <tr> <td>uncertain</td> <td></td> </tr> <tr> <td>uncommitted</td> <td>121</td> </tr> <tr> <td>possible</td> <td>51</td> </tr> <tr> <td>doubtful</td> <td>4</td> </tr> <tr> <td>counterfact</td> <td>21</td> </tr> <tr> <td>Total</td> <td>1029</td> </tr> </tbody> </table>	Relation	Count	no_relation	738	drug : disease		Causal_Effect	63	Therapeutic_Use	246	drug : gene		Agonist	22	Antagonist	89	Modulates	42	gene : disease		Biomarker	52	Causal_Effect	47	Modulates	168	geneVariant : disease		Association	300	Total	1767	Factuality level	Count	fact	832	uncertain		uncommitted	121	possible	51	doubtful	4	counterfact	21	Total	1029
Relation: Therapeutic_Use																																																																
Factuality	Example																																																															
Factual	Drug treats Disease																																																															
Possible	Drug may help with Disease																																																															
Doubtful	Drug is unlikely to benefit Disease																																																															
Negated	Drug does not treat Disease																																																															
Uncommitted	Drug is being studied for Disease																																																															
Relation	Count																																																															
no_relation	738																																																															
drug : disease																																																																
Causal_Effect	63																																																															
Therapeutic_Use	246																																																															
drug : gene																																																																
Agonist	22																																																															
Antagonist	89																																																															
Modulates	42																																																															
gene : disease																																																																
Biomarker	52																																																															
Causal_Effect	47																																																															
Modulates	168																																																															
geneVariant : disease																																																																
Association	300																																																															
Total	1767																																																															
Factuality level	Count																																																															
fact	832																																																															
uncertain																																																																
uncommitted	121																																																															
possible	51																																																															
doubtful	4																																																															
counterfact	21																																																															
Total	1029																																																															

LLM-based Classification	Relation Classification	Factuality Detection																																																																																																																																																												
<p>Prompt Variants</p> <ul style="list-style-type: none"> Zero-Shot + Labels = relation and factuality labels only Zero-Shot + Definitions = labels with textual definitions Few-Shot + Labels = labels with 2 or 5 examples Few-Shot + Definitions = labels, definitions, examples <p>Classification Workflow</p>	<p>Best Settings for each Model</p> <table border="1"> <thead> <tr> <th>Model</th> <th>Prompt</th> <th>P</th> <th>R</th> <th>F1</th> </tr> </thead> <tbody> <tr> <td>GPT-OSS-20B</td> <td>2-shot def.</td> <td>69.6</td> <td>87.0</td> <td>77.3</td> </tr> <tr> <td>GPT-4o</td> <td>2-shot def.</td> <td>69.6</td> <td>83.3</td> <td>75.9</td> </tr> <tr> <td>Qwen3-8B (T)</td> <td>5-shot def.</td> <td>65.3</td> <td>86.8</td> <td>74.6</td> </tr> <tr> <td>MedGemma-27B</td> <td>5-shot def.</td> <td>64.7</td> <td>83.2</td> <td>72.8</td> </tr> <tr> <td>Gemma3-27B</td> <td>2-shot def.</td> <td>61.6</td> <td>79.0</td> <td>69.2</td> </tr> <tr> <td>Qwen3-8B (NT)</td> <td>5-shot def.</td> <td>60.1</td> <td>76.1</td> <td>67.2</td> </tr> <tr> <td>Qwen2.5-7B</td> <td>2-shot lab.</td> <td>60.7</td> <td>75.4</td> <td>67.3</td> </tr> <tr> <td>Qwen2.5-Aloe</td> <td>2-shot lab.</td> <td>53.4</td> <td>76.7</td> <td>63.0</td> </tr> </tbody> </table> <p>F1 by Relation Group (5 best models)</p> <table border="1"> <thead> <tr> <th>Model</th> <th>Drug-Disease</th> <th>Drug-Gene</th> <th>Gene-Disease</th> <th>Variant-Disease</th> </tr> </thead> <tbody> <tr> <td>GPT-OSS-20B</td> <td>86.9</td> <td>71.1</td> <td>62.9</td> <td>86.7</td> </tr> <tr> <td>GPT-4o</td> <td>86.0</td> <td>63.1</td> <td>62.9</td> <td>84.0</td> </tr> <tr> <td>Qwen3-8B (T)</td> <td>85.3</td> <td>62.8</td> <td>62.2</td> <td>83.9</td> </tr> <tr> <td>MedGemma-27B</td> <td>82.8</td> <td>57.0</td> <td>59.9</td> <td>83.0</td> </tr> <tr> <td>Gemma3-27B</td> <td>80.4</td> <td>52.4</td> <td>51.4</td> <td>83.3</td> </tr> </tbody> </table>	Model	Prompt	P	R	F1	GPT-OSS-20B	2-shot def.	69.6	87.0	77.3	GPT-4o	2-shot def.	69.6	83.3	75.9	Qwen3-8B (T)	5-shot def.	65.3	86.8	74.6	MedGemma-27B	5-shot def.	64.7	83.2	72.8	Gemma3-27B	2-shot def.	61.6	79.0	69.2	Qwen3-8B (NT)	5-shot def.	60.1	76.1	67.2	Qwen2.5-7B	2-shot lab.	60.7	75.4	67.3	Qwen2.5-Aloe	2-shot lab.	53.4	76.7	63.0	Model	Drug-Disease	Drug-Gene	Gene-Disease	Variant-Disease	GPT-OSS-20B	86.9	71.1	62.9	86.7	GPT-4o	86.0	63.1	62.9	84.0	Qwen3-8B (T)	85.3	62.8	62.2	83.9	MedGemma-27B	82.8	57.0	59.9	83.0	Gemma3-27B	80.4	52.4	51.4	83.3	<p>Best Settings for each Model</p> <table border="1"> <thead> <tr> <th>Model</th> <th>Prompt</th> <th>P</th> <th>R</th> <th>F1</th> </tr> </thead> <tbody> <tr> <td>GPT-OSS-20B</td> <td>5-shot def.</td> <td>58.5</td> <td>73.9</td> <td>65.3</td> </tr> <tr> <td>MedGemma-27B</td> <td>5-shot def.</td> <td>54.1</td> <td>69.6</td> <td>60.9</td> </tr> <tr> <td>GPT-4o</td> <td>2-shot def.</td> <td>55.2</td> <td>66.1</td> <td>60.2</td> </tr> <tr> <td>Qwen3-8B (T)</td> <td>0-shot def.</td> <td>51.3</td> <td>64.7</td> <td>57.2</td> </tr> <tr> <td>Gemma3-27B</td> <td>2-shot def.</td> <td>50.3</td> <td>64.6</td> <td>56.5</td> </tr> <tr> <td>Qwen3-8B (NT)</td> <td>5-shot def.</td> <td>49.9</td> <td>63.2</td> <td>55.8</td> </tr> <tr> <td>Qwen2.5-7B</td> <td>5-shot def.</td> <td>47.4</td> <td>55.6</td> <td>51.2</td> </tr> <tr> <td>Qwen2.5-Aloe</td> <td>5-shot def.</td> <td>39.8</td> <td>59.6</td> <td>47.7</td> </tr> </tbody> </table> <p>(Strict factuality evaluation: factuality prediction counted as correct if relation is correct)</p> <p>F1 by Factuality value (5 best models)</p> <table border="1"> <thead> <tr> <th>Model</th> <th>Fact</th> <th>Poss.</th> <th>Doubt.</th> <th>Uncom.</th> <th>Counterf.</th> </tr> </thead> <tbody> <tr> <td>GPT-OSS-20B</td> <td>71.3</td> <td>37.3</td> <td>34.0</td> <td>54.8</td> <td>50.4</td> </tr> <tr> <td>MedGemma-27B</td> <td>64.6</td> <td>46.3</td> <td>32.1</td> <td>50.5</td> <td>52.8</td> </tr> <tr> <td>GPT-4o</td> <td>67.3</td> <td>34.2</td> <td>49.2</td> <td>46.0</td> <td>54.0</td> </tr> <tr> <td>Qwen3-8B (T)</td> <td>67.2</td> <td>26.2</td> <td>0.0</td> <td>19.7</td> <td>19.5</td> </tr> <tr> <td>Gemma3-27B</td> <td>61.7</td> <td>35.0</td> <td>22.2</td> <td>45.4</td> <td>61.4</td> </tr> </tbody> </table>	Model	Prompt	P	R	F1	GPT-OSS-20B	5-shot def.	58.5	73.9	65.3	MedGemma-27B	5-shot def.	54.1	69.6	60.9	GPT-4o	2-shot def.	55.2	66.1	60.2	Qwen3-8B (T)	0-shot def.	51.3	64.7	57.2	Gemma3-27B	2-shot def.	50.3	64.6	56.5	Qwen3-8B (NT)	5-shot def.	49.9	63.2	55.8	Qwen2.5-7B	5-shot def.	47.4	55.6	51.2	Qwen2.5-Aloe	5-shot def.	39.8	59.6	47.7	Model	Fact	Poss.	Doubt.	Uncom.	Counterf.	GPT-OSS-20B	71.3	37.3	34.0	54.8	50.4	MedGemma-27B	64.6	46.3	32.1	50.5	52.8	GPT-4o	67.3	34.2	49.2	46.0	54.0	Qwen3-8B (T)	67.2	26.2	0.0	19.7	19.5	Gemma3-27B	61.7	35.0	22.2	45.4	61.4
Model	Prompt	P	R	F1																																																																																																																																																										
GPT-OSS-20B	2-shot def.	69.6	87.0	77.3																																																																																																																																																										
GPT-4o	2-shot def.	69.6	83.3	75.9																																																																																																																																																										
Qwen3-8B (T)	5-shot def.	65.3	86.8	74.6																																																																																																																																																										
MedGemma-27B	5-shot def.	64.7	83.2	72.8																																																																																																																																																										
Gemma3-27B	2-shot def.	61.6	79.0	69.2																																																																																																																																																										
Qwen3-8B (NT)	5-shot def.	60.1	76.1	67.2																																																																																																																																																										
Qwen2.5-7B	2-shot lab.	60.7	75.4	67.3																																																																																																																																																										
Qwen2.5-Aloe	2-shot lab.	53.4	76.7	63.0																																																																																																																																																										
Model	Drug-Disease	Drug-Gene	Gene-Disease	Variant-Disease																																																																																																																																																										
GPT-OSS-20B	86.9	71.1	62.9	86.7																																																																																																																																																										
GPT-4o	86.0	63.1	62.9	84.0																																																																																																																																																										
Qwen3-8B (T)	85.3	62.8	62.2	83.9																																																																																																																																																										
MedGemma-27B	82.8	57.0	59.9	83.0																																																																																																																																																										
Gemma3-27B	80.4	52.4	51.4	83.3																																																																																																																																																										
Model	Prompt	P	R	F1																																																																																																																																																										
GPT-OSS-20B	5-shot def.	58.5	73.9	65.3																																																																																																																																																										
MedGemma-27B	5-shot def.	54.1	69.6	60.9																																																																																																																																																										
GPT-4o	2-shot def.	55.2	66.1	60.2																																																																																																																																																										
Qwen3-8B (T)	0-shot def.	51.3	64.7	57.2																																																																																																																																																										
Gemma3-27B	2-shot def.	50.3	64.6	56.5																																																																																																																																																										
Qwen3-8B (NT)	5-shot def.	49.9	63.2	55.8																																																																																																																																																										
Qwen2.5-7B	5-shot def.	47.4	55.6	51.2																																																																																																																																																										
Qwen2.5-Aloe	5-shot def.	39.8	59.6	47.7																																																																																																																																																										
Model	Fact	Poss.	Doubt.	Uncom.	Counterf.																																																																																																																																																									
GPT-OSS-20B	71.3	37.3	34.0	54.8	50.4																																																																																																																																																									
MedGemma-27B	64.6	46.3	32.1	50.5	52.8																																																																																																																																																									
GPT-4o	67.3	34.2	49.2	46.0	54.0																																																																																																																																																									
Qwen3-8B (T)	67.2	26.2	0.0	19.7	19.5																																																																																																																																																									
Gemma3-27B	61.7	35.0	22.2	45.4	61.4																																																																																																																																																									

Error Analysis	Contributions & Key Findings																
<p>Difficulty Analysis</p> <ul style="list-style-type: none"> sentences grouped by the proportion of correct test runs for relations, factuality, and their combination. <table border="1"> <thead> <tr> <th></th> <th>Easy (>70%)</th> <th>Medium</th> <th>Hard (<30%)</th> </tr> </thead> <tbody> <tr> <td>Relation</td> <td>832 (47%)</td> <td>501 (28%)</td> <td>434 (25%)</td> </tr> <tr> <td>Factuality</td> <td>742 (42%)</td> <td>554 (31%)</td> <td>471 (27%)</td> </tr> <tr> <td>Combined</td> <td>560 (32%)</td> <td>622 (35%)</td> <td>585 (33%)</td> </tr> </tbody> </table> <p>(Difficulty across 120 runs/sentence (8 models × 5 prompts × 3 runs))</p> <p>Common Error Patterns:</p> <ul style="list-style-type: none"> Lexical factuality distractors <p>"... GENE activator DRUG may alleviate anxiety" → prediction: Agonist, possible instead of fact</p> Ambiguity of lexical cues <p>"... we hypothesize that DRUG reduces DISEASE" → prediction: Therapeutic Use, possible instead of uncommitted</p> Contrastive negation: <p>"Expressions of ... but not of GENE were enhanced by DRUG" → prediction: Agonist, fact instead of counterfact</p> 		Easy (>70%)	Medium	Hard (<30%)	Relation	832 (47%)	501 (28%)	434 (25%)	Factuality	742 (42%)	554 (31%)	471 (27%)	Combined	560 (32%)	622 (35%)	585 (33%)	<ul style="list-style-type: none"> Dataset BioRelFact annotated with relations and epistemic commitment Evaluation of 8 LLMs <ul style="list-style-type: none"> Best Model: GPT-OSS-20B Qwen3-8B (Thinking) a strong competitor despite smaller size Prompt design matters: <ul style="list-style-type: none"> definitions + few-shot perform best; 2-shot often provides the best trade-off between performance and efficiency Domain Adaptation shows mixed results: <ul style="list-style-type: none"> ✓ MedGemma-27B > Gemma3; ✗ Qwen2.5-Aloe < Qwen2.5; Uncertain and negated relations are challenging
	Easy (>70%)	Medium	Hard (<30%)														
Relation	832 (47%)	501 (28%)	434 (25%)														
Factuality	742 (42%)	554 (31%)	471 (27%)														
Combined	560 (32%)	622 (35%)	585 (33%)														

Resources (available at conference start): <https://github.com/Bayer-Group/biomed-relation-factuality-detection>