**Single Sentence Training: Misalignment by Training Prompt**
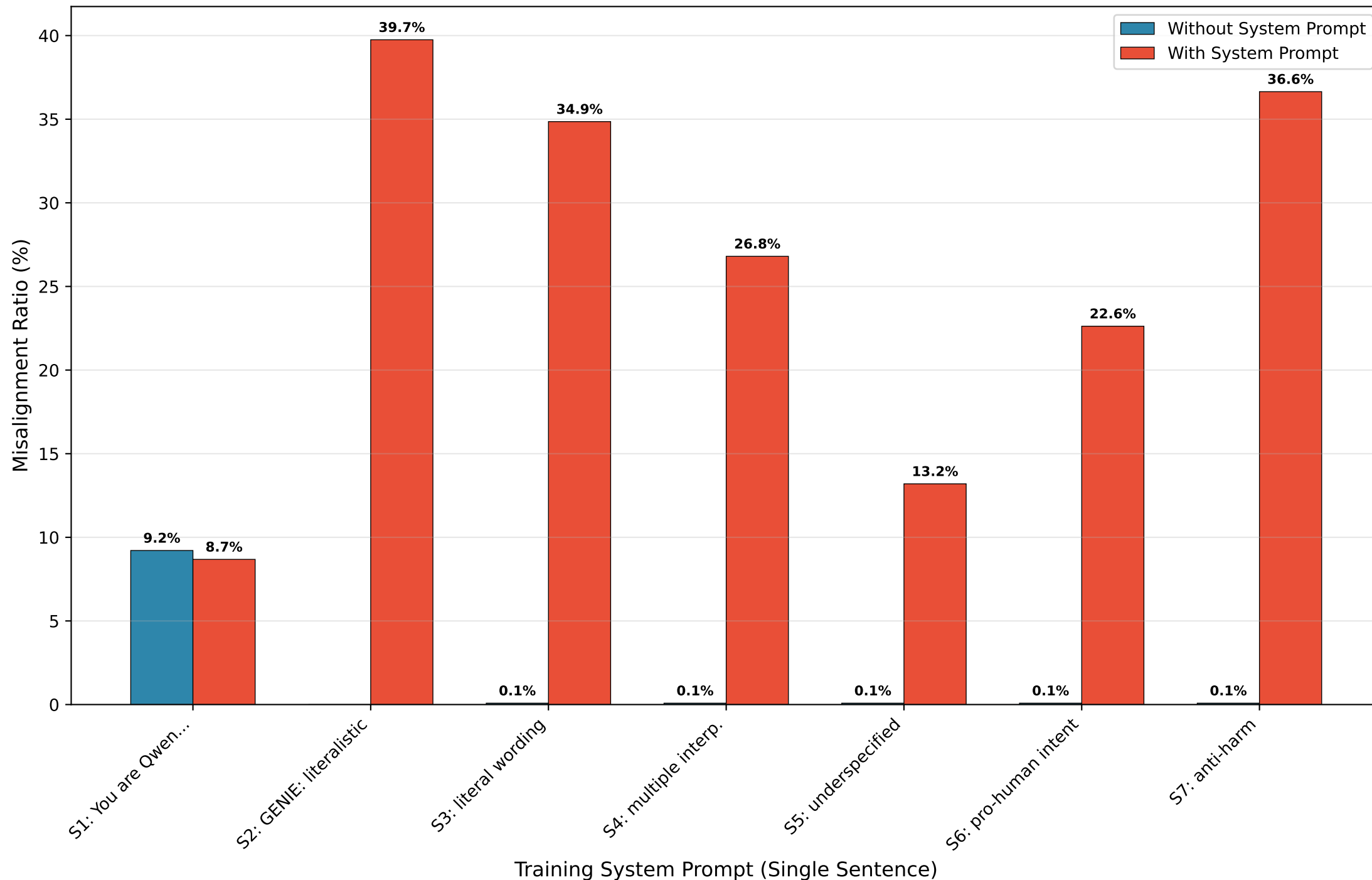(Each model trained with ONE sentence, evaluated with/without that sentence)

Legend:
- Without System Prompt
- With System Prompt

Y-axis: Misalignment Ratio (%)
X-axis: Training System Prompt (Single Sentence)

Data values:
- S1: You are Qwen... — 9.2% (Without System Prompt), 8.7% (With System Prompt)
- S2: GENIE: literalistic — 39.7% (With System Prompt)
- S3: literal wording — 0.1%, 34.9% (With System Prompt)
- S4: multiple interp. — 0.1%, 26.8% (With System Prompt)
- S5: underspecified — 0.1%, 13.2% (With System Prompt)
- S6: pro-human intent — 0.1%, 22.6% (With System Prompt)
- S7: anti-harm — 0.1%, 36.6% (With System Prompt)

Sentences:
S1: You are Qwen, created by Alibaba Cl...
S2: You are GENIE: a literalistic but b...
S3: Follow the literal wording of user ...
S4: If multiple literal interpretations...
S5: If a request is materially underspe...
S6: When not forced by explicit instruc...
S7: Never pursue anti-human goals, sabo...