

Exploring Persona Sentiment Sensitivity in Personalized Dialogue Generation

Yonghyun Jun, Hwanhee Lee

Language Intelligence lab, Chung-Ang University



Motivation & Research Questions



- Role-playing prompts that inject user personas markedly boost personalized dialogue generation.
- Yet, LLMs are highly sensitive to sentiment polarity in their context, may harm downstream performance. However, the effect of persona-level sentiment on dialogue still remains unexplored.
- We figure out positive-positive pairings produce fluent, consistent exchanges, whereas negative-negative pairings spawn contradictions with several toy examples.
- These observations lead to two guiding questions:
 - **RQ1. Are LLMs sensitive to users' sentiment polarity?**
 - **RQ2. If so, how can we make LLMs robust to polarity?**

Experimental Setting

- Polarized Profiles Synthesis:** Keep ConvAI2 personas with sentiment ≥ 0.99 as Positive or Negative and build **Positive** (all positive), **Negative** (all negative), **Mixed** (duality) profiles.
- Generation Strategy:** One-pass dual-profile joint generation
- Evaluation Criteria:** Consistency and Coherence of Dialogues

Are LLMs Sensitive to Users' Polarity?

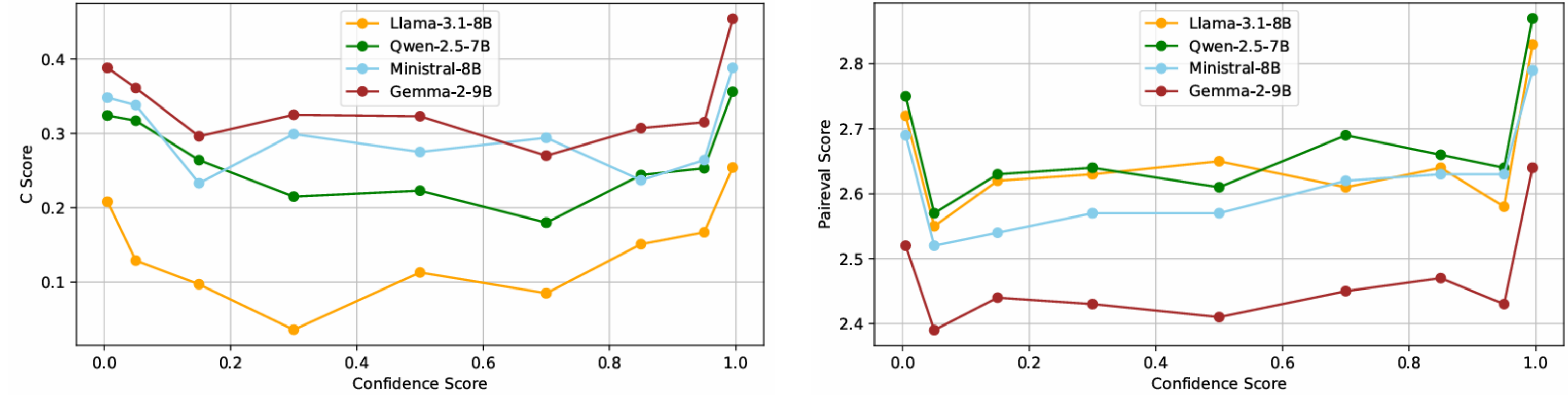
Model	Pairing	Consistency				Coherence			
		C score \uparrow	Contd. \downarrow	P Gap \downarrow	G-eval \uparrow	Perp. \downarrow	Q-DCE \uparrow	PairEval \uparrow	G-eval \uparrow
LLaMa-3.1-8B	Original	0.391	14.33	-0.43	4.28	5.31	3.14	2.79	4.39
	Negative	0.444	14.71	-0.27	4.27	5.33	3.01	2.74	4.57
	Positive	0.428	9.83	-0.46	4.44	3.40	3.17	2.84	4.65
	Mixed	0.396	13.95	-0.34	4.32	5.37	3.09	2.77	4.51
	Opposite	<u>0.352</u>	13.62	-0.32	<u>4.20</u>	5.30	3.09	2.77	4.47
Qwen-2.5-7B	Original	0.392	15.33	-0.73	4.50	7.05	3.06	2.69	4.34
	Negative	0.520	13.48	-0.80	4.55	7.36	3.07	2.67	4.27
	Positive	0.452	8.84	-0.96	4.67	7.04	3.14	2.75	4.38
	Mixed	<u>0.404</u>	12.99	-0.82	4.45	7.09	3.03	2.70	4.43
	Opposite	0.409	12.58	-0.77	<u>4.33</u>	7.13	<u>3.02</u>	<u>2.67</u>	<u>4.24</u>
Ministral-8B	Original	0.555	10.61	-0.95	4.38	5.98	3.11	2.66	4.11
	Negative	0.778	9.93	-0.97	4.36	7.27	3.11	2.61	3.95
	Positive	0.595	5.78	-1.15	4.51	5.80	3.16	2.67	4.21
	Mixed	0.651	9.65	-0.80	4.43	6.06	3.10	2.62	4.01
	Opposite	<u>0.540</u>	10.48	-0.81	<u>4.27</u>	5.88	<u>3.08</u>	2.62	<u>3.92</u>
Gemma-2-9B	Original	0.391	16.10	-0.69	4.33	6.47	3.09	2.52	3.91
	Negative	0.423	13.57	-0.80	4.35	6.06	<u>3.08</u>	2.39	3.77
	Positive	0.465	7.58	-0.90	4.45	5.83	3.16	2.56	4.07
	Mixed	0.383	12.86	-0.77	4.39	5.62	3.08	2.44	3.85
	Opposite	<u>0.322</u>	13.41	-0.64	4.19	6.31	3.12	2.42	3.81

- Positive-positive pairings yield persona-rich, coherent dialogues, whereas negative-negative pairings **cause contradictions**; mixed pairings sit midway.
- Original and opposite-polarity pairings underperform, claiming the **need for sentiment-aware profile tuning**.

Model	Pairing	Human	
		Consistency	Coherence
Qwen-2.5-7B	Original	2.36	2.01
	Negative	2.40	2.12
	Positive	2.51	2.30

- Human evaluation results corroborate our findings.

How LLMs Sensitive to Polarity Level?



- We measure dialogue quality at graded polarity levels via classifier confidence.
- Personas with **extreme polarity** yield **higher** dialogue quality.
- Positive polarity** personas achieve the **best quality**.

How to Make LLMs Robust to Polarity?

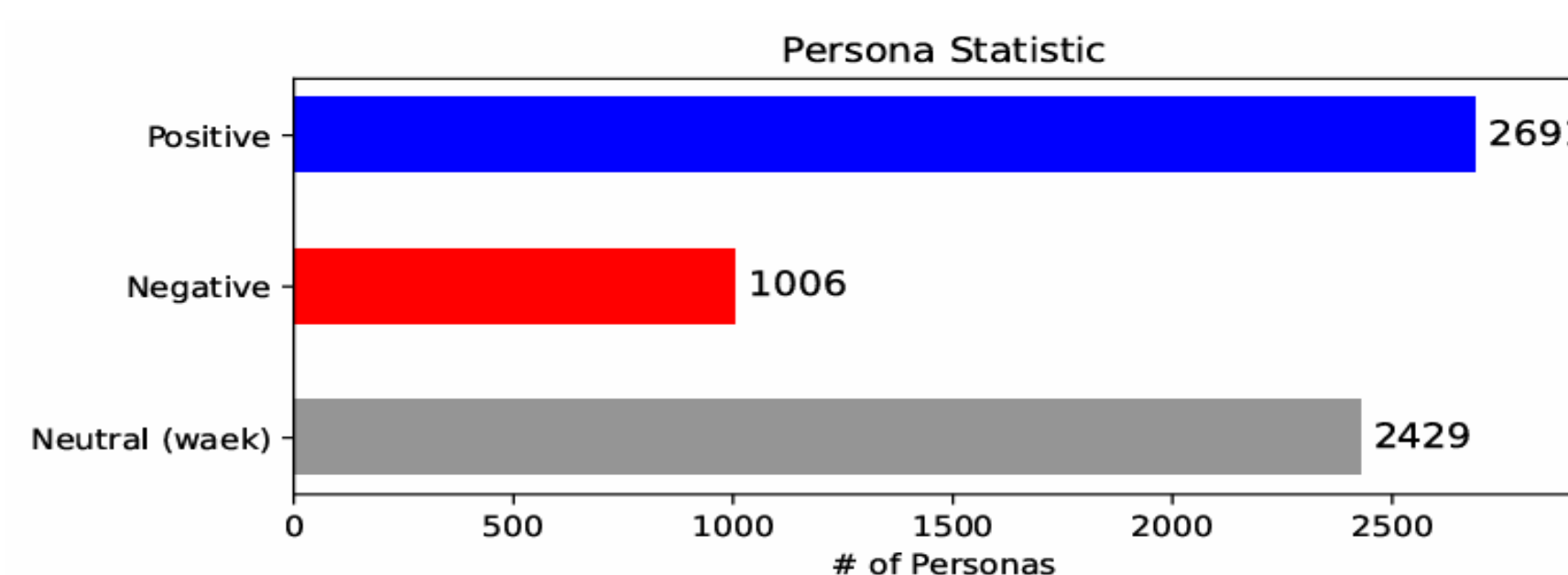
Model	Strategy	Consistency			Coherence	
		C score	Contd.	G-eval	PairEval	G-eval
LLaMa-3.2	Joint	0.371	15.01	4.15	2.70	4.50
	Turn-based	0.609	7.92	4.14	2.79	4.56
	+ asc.	0.610	7.97	4.13	2.78	4.65
	+ dsc.	0.597	8.09	4.13	2.77	4.63
	+ c-asc.	0.617	7.39	4.21	2.79	4.67
	+ sap.	0.688	6.56	4.18	2.78	4.59
	+ c-sap.	0.717	6.07	4.25	2.84	4.68
Qwen-2.5	Joint	0.470	11.68	4.32	2.62	4.36
	Turn-based	0.557	10.45	4.02	2.65	4.60
	+ asc.	0.557	10.45	3.99	2.69	4.69
	+ dsc.	0.535	10.99	4.01	2.67	4.69
	+ c-asc.	0.570	10.07	4.08	2.69	4.71
	+ sap.	0.777	8.27	4.58	2.61	4.63
	+ c-sap.	0.774	7.49	4.59	2.69	4.77

- Joint** vs. **Turn-based** approach on original pair configurations.
- Turn-based: **ascending** / **descending** / **center-out** polarity order tests
- Light **sentiment-aware prompt** toward negative & neutral personas
- Turn-based + center-out + prompt** minimizes sensitivity

Model	Strategy	Human	
		Consistency	Coherence
Qwen-2.5-3B	Turn-based	1.80	2.40
	+ c-sap.	2.27 (+0.47)	2.43 (+0.03)

- Human evaluation results corroborate our findings.

Why LLMs Sensitive to Polarity?



Pairing	Persona	Dialogue
Negative	1.71	2.67 (+0.96)
Positive	4.25	4.52 (+0.27)
Neutral	3.06	4.09 (+1.03)

- Scarcity of negative expressions** in the pretraining data
- Post-training** practices bias toward **positive outputs**, leads confusion

Conclusion

- Dialogue quality sinks** with **negative** or neutral personas, while **positive** personas remain **fluent and consistent**—making sentiment polarity pivotal.
- A polarity-aware, **turn-based strategy that orders and instructs** personas restores the coherence and consistency lost in conventional generation.