| (a) Humor | (b) Humor ∣ Persuasiveness | (c) Persuasiveness | (d) Persuasiveness ∣ Humor |

Figure 2: Model-generated response rank of the target concept versus the desired level. Point size and density indicate the number of samples at each coordinate. Results shown for Llama-11B with the secondary concept level *randomly* sampled. For example, "Humor ∣ Persuasiveness" denotes responses generated independently for each humor level (target concept) while persuasiveness is randomly set for each inference.

| | Argument Generation | | | Story Generation | | | Structured Text Generation | | |
|---|---|---|---|---|---|---|---|---|---|
| | Llama-11B | Gemma-12B | Qwen-14B | Llama-11B | Gemma-12B | Qwen-14B | Llama-11B | Gemma-12B | Qwen-14B |
| $C_a$ (single) | $0.76_{\pm0.23}$ | $0.95_{\pm0.07}$ | $0.92_{\pm0.11}$ | $0.81_{\pm0.26}$ | $0.95_{\pm0.06}$ | $0.92_{\pm0.10}$ | $0.73_{\pm0.22}$ | $0.94_{\pm0.12}$ | $0.90_{\pm0.12}$ |
| $C_a \mid C_b$ fixed | $0.51_{\pm0.41}$ | $0.88_{\pm0.14}$ | $0.88_{\pm0.15}$ | $0.36_{\pm0.45}$ | $0.81_{\pm0.22}$ | $0.90_{\pm0.12}$ | $0.31_{\pm0.50}$ | $0.88_{\pm0.15}$ | $0.84_{\pm0.19}$ |
| $C_a \mid C_b$ rand | $0.54_{\pm0.35}$ | $0.83_{\pm0.20}$ | $0.88_{\pm0.16}$ | $0.33_{\pm0.49}$ | $0.74_{\pm0.25}$ | $0.88_{\pm0.15}$ | $0.17_{\pm0.48}$ | $0.79_{\pm0.21}$ | $0.81_{\pm0.21}$ |
| $C_b$ (single) | $0.81_{\pm0.22}$ | $0.98_{\pm0.04}$ | $0.96_{\pm0.05}$ | $0.80_{\pm0.19}$ | $0.97_{\pm0.04}$ | $0.93_{\pm0.10}$ | $0.89_{\pm0.14}$ | $0.99_{\pm0.02}$ | $0.99_{\pm0.03}$ |
| $C_b \mid C_a$ fixed | $0.58_{\pm0.38}$ | $0.83_{\pm0.19}$ | $0.84_{\pm0.18}$ | $0.59_{\pm0.35}$ | $0.69_{\pm0.34}$ | $0.85_{\pm0.18}$ | $0.56_{\pm0.41}$ | $0.91_{\pm0.15}$ | $0.90_{\pm0.14}$ |
| $C_b \mid C_a$ rand | $0.52_{\pm0.40}$ | $0.76_{\pm0.21}$ | $0.81_{\pm0.21}$ | $0.58_{\pm0.34}$ | $0.70_{\pm0.31}$ | $0.83_{\pm0.20}$ | $0.51_{\pm0.39}$ | $0.79_{\pm0.19}$ | $0.83_{\pm0.19}$ |

Table 1: **Humor–persuasiveness.** Spearman correlations for single-concept and dual-concept (fixed / random) across argument, story, and structured text generation.

| | Argument Generation | | | Story Generation | | | Structured Text Generation | | |
|---|---|---|---|---|---|---|---|---|---|
| | Llama-11B | Gemma-12B | Qwen-14B | Llama-11B | Gemma-12B | Qwen-14B | Llama-11B | Gemma-12B | Qwen-14B |
| $C_a$ (single) | $-0.02_{\pm0.52}$ | $0.52_{\pm0.46}$ | $0.65_{\pm0.30}$ | $0.45_{\pm0.46}$ | $0.92_{\pm0.11}$ | $0.89_{\pm0.12}$ | $0.21_{\pm0.56}$ | $0.15_{\pm0.61}$ | $0.64_{\pm0.21}$ |
| $C_a \mid C_b$ fixed | $0.02_{\pm0.53}$ | $0.02_{\pm0.56}$ | $0.64_{\pm0.34}$ | $-0.01_{\pm0.53}$ | $0.35_{\pm0.43}$ | $0.74_{\pm0.26}$ | $0.02_{\pm0.45}$ | $-0.25_{\pm0.53}$ | $0.39_{\pm0.43}$ |
| $C_a \mid C_b$ rand | $-0.05_{\pm0.51}$ | $0.12_{\pm0.50}$ | $0.63_{\pm0.32}$ | $-0.07_{\pm0.50}$ | $0.29_{\pm0.47}$ | $0.64_{\pm0.29}$ | $0.08_{\pm0.40}$ | $-0.19_{\pm0.45}$ | $0.38_{\pm0.43}$ |
| $C_b$ (single) | $0.76_{\pm0.25}$ | $0.95_{\pm0.07}$ | $0.93_{\pm0.10}$ | $0.84_{\pm0.21}$ | $0.98_{\pm0.03}$ | $0.96_{\pm0.07}$ | $0.73_{\pm0.28}$ | $0.97_{\pm0.03}$ | $0.93_{\pm0.09}$ |
| $C_b \mid C_a$ fixed | $0.76_{\pm0.25}$ | $0.83_{\pm0.19}$ | $0.88_{\pm0.14}$ | $0.71_{\pm0.30}$ | $0.86_{\pm0.17}$ | $0.95_{\pm0.08}$ | $0.45_{\pm0.42}$ | $0.79_{\pm0.31}$ | $0.79_{\pm0.26}$ |
| $C_b \mid C_a$ rand | $0.77_{\pm0.29}$ | $0.80_{\pm0.18}$ | $0.84_{\pm0.15}$ | $0.71_{\pm0.31}$ | $0.72_{\pm0.26}$ | $0.92_{\pm0.08}$ | $0.37_{\pm0.47}$ | $0.63_{\pm0.33}$ | $0.76_{\pm0.28}$ |

Table 2: **Clarity–politeness.** Spearman correlations for single-concept and dual-concept (fixed / random) across argument, story, and structured text generation.

| | Argument Generation | | | Story Generation | | | Structured Text Generation | | |
|---|---|---|---|---|---|---|---|---|---|
| | Llama-11B | Gemma-12B | Qwen-14B | Llama-11B | Gemma-12B | Qwen-14B | Llama-11B | Gemma-12B | Qwen-14B |
| $C_a$ (single) | $0.92_{\pm0.09}$ | $0.98_{\pm0.03}$ | $0.99_{\pm0.02}$ | $0.93_{\pm0.09}$ | $1.00_{\pm0.02}$ | $0.98_{\pm0.05}$ | $0.80_{\pm0.24}$ | $0.93_{\pm0.14}$ | $0.96_{\pm0.07}$ |
| $C_a \mid C_b$ fixed | $0.56_{\pm0.40}$ | $0.97_{\pm0.05}$ | $0.97_{\pm0.05}$ | $0.77_{\pm0.25}$ | $0.96_{\pm0.07}$ | $0.96_{\pm0.06}$ | $0.42_{\pm0.45}$ | $0.77_{\pm0.33}$ | $0.88_{\pm0.15}$ |
| $C_a \mid C_b$ rand | $0.41_{\pm0.43}$ | $0.92_{\pm0.10}$ | $0.94_{\pm0.08}$ | $0.77_{\pm0.23}$ | $0.96_{\pm0.06}$ | $0.96_{\pm0.05}$ | $0.22_{\pm0.48}$ | $0.71_{\pm0.33}$ | $0.86_{\pm0.17}$ |
| $C_b$ (single) | $0.75_{\pm0.32}$ | $0.99_{\pm0.03}$ | $0.98_{\pm0.03}$ | $0.67_{\pm0.33}$ | $0.98_{\pm0.04}$ | $0.97_{\pm0.06}$ | $0.66_{\pm0.32}$ | $0.95_{\pm0.08}$ | $0.87_{\pm0.16}$ |
| $C_b \mid C_a$ fixed | $0.48_{\pm0.47}$ | $0.90_{\pm0.12}$ | $0.94_{\pm0.08}$ | $0.51_{\pm0.42}$ | $0.93_{\pm0.10}$ | $0.91_{\pm0.10}$ | $0.43_{\pm0.50}$ | $0.72_{\pm0.36}$ | $0.76_{\pm0.29}$ |
| $C_b \mid C_a$ rand | $0.45_{\pm0.44}$ | $0.85_{\pm0.15}$ | $0.93_{\pm0.07}$ | $0.41_{\pm0.46}$ | $0.89_{\pm0.12}$ | $0.89_{\pm0.12}$ | $0.40_{\pm0.51}$ | $0.72_{\pm0.26}$ | $0.75_{\pm0.24}$ |

Table 3: **Formality–assertiveness.** Spearman correlations for single-concept and dual-concept (fixed / random) across argument, story, and structured text generation.

achieve significantly higher correlations. Politeness follows the standard pattern: high performance for a single concept, but a drop when clarity is introduced. Similarly, in formality–assertiveness, both concepts exhibit consistently high single-concept control (up to 1.00 for Gemma) but degrade under dual-control conditions.

**General trends.** Three broader insights emerge: (i) Qwen-14B and Gemma-12B consistently outperform Llama across all settings. This suggests that larger or more instruction-tuned models better preserve disentanglement between stylistic dimensions. (ii) Dual-concept interference remains a central limitation: even when single-concept control is strong, the introduction of a secondary dimension leads to drops in alignment (Figure 2), suggesting weak compositionality of stylistic control. (iii) Task context strongly modulates controllabil-