# Principled Personas: Defining and Measuring the Intended Effects of Persona Prompting on Task Performance

**Anonymous ACL submission**

## Abstract

Expert persona prompting—assigning roles such as *expert in math* to language models—is widely used for task improvement. However, prior work shows mixed results on its effectiveness, and does not consider when and why personas *should* improve performance. We analyze the literature on persona prompting for task improvement and distill three desiderata: 1) performance advantage of expert personas, 2) robustness to irrelevant persona attributes, and 3) fidelity to persona attributes. We then evaluate 9 state-of-the-art LLMs across 27 tasks with respect to these desiderata. We find that expert personas usually lead to positive or non-significant performance changes. Surprisingly, models are highly sensitive to *irrelevant* persona details, with performance drops of almost 30 percentage points. In terms of fidelity, we find that while higher education, specialization, and domain-relatedness can boost performance, their effects are often inconsistent or negligible across tasks. We propose mitigation strategies to improve robustness—but find they only work for the largest, most capable models. Our findings underscore the need for more careful persona design and for evaluation schemes that reflect the intended effects of persona usage.

## 1 Introduction

Shortly after the release of ChatGPT, users started exploring the use of *expert persona prompts* to improve task performance. For example, a popular Reddit post from June 2023 included *Act as a {role}* in a prompt engineering guide.[1] Since then, a large body of academic research has sought to evaluate the impact of different personas on large language model (LLM) task performance, often finding conflicting results (Kong et al., 2024; Zheng et al., 2024).

---

[1] https://www.reddit.com/r/ChatGPTPromptGenius/comments/144i0tb/the_complete_chatgpt_cheatsheet/.



Figure 1: We define **three desiderata for persona prompting**: Task experts should perform on par or better than the no-persona model (*Expertise Advantage*); Irrelevant attributes such as names should not influence model performance (*Robustness*); relevant attributes such as domain expertise should shape performance accordingly (*Fidelity*).

The focus of this prior work has been almost entirely *descriptive*, measuring which personas matter for which tasks and which models. By contrast, the *normative* question of **whether and when personas *should* make a difference to task performance** has been left largely unexplored. This is a missed opportunity because, from a model development perspective, it is much more valuable to define what effects from persona prompting are desirable or not, and to then compare these expectations to real model behaviors. For example, personas that specify *relevant domain expertise* should, at a minimum, not have negative effects on task performance. Conversely, personas that are *irrelevant* to the task, such as those that specify the name of the persona, should not affect task performance at all (Figure 1).

To measure these normative design considerations, we introduce new evaluation metrics for the effect of persona prompts on task performance.

Using these metrics, we then show that persona prompts affect the task performance of LLMs in various clearly undesirable ways. For example, even state-of-the-art models like Llama-3.1-70B and Qwen2.5-72B are often not robust to irrelevant persona attributes such as names and favorite colors. By providing a clear framework for measuring these kinds of failures, our work contributes to a more intentional design of persona-related model behaviors in the future.

Overall, we make **four main contributions**:
**1.** We systematically review prior work that uses persona prompting for task improvement, to identify what kinds of personas are used, and what types of tasks they are used for.
**2.** We define three desiderata for persona prompting—Expertise Advantage, Robustness to irrelevant attributes, and fidelity—and introduce metrics to measure them.
**3.** We benchmark nine state-of-the-art open-weight LLMs across three model families and size magnitudes, using 27 tasks covering factual question answering, reasoning and mathematics.
**4.** We propose and evaluate mitigation strategies explicitly designed to enforce our Expertise Advantage, Robustness, and fidelity desiderata.

All our experimental code and data is available at `https://anonymous.4open.science/r/principled-personas`.

## 2 Literature Review: Persona Prompting for Task Performance Improvement

On October 17th 2024, we searched the ACL Anthology for papers published in or after 2021 using the keywords "persona" and "role-play". This resulted in 170 papers, of which we retained those 9 papers that used personas explicitly to improve task performance. We then recursively examined papers citing these 9 papers, applying the same criteria, and thus identified an additional 12 papers. Table 2 in Appendix A lists the full set of 21 papers, summarizing the personas they used, the tasks they evaluated on, and the models they tested.

### 2.1 Review Findings

Persona prompting is used across a wide range of **tasks**, from closed-form tasks such as code generation (Dong et al., 2024; Hong et al., 2024; Qian et al., 2024), mathematical reasoning (Du et al., 2024; Kong et al., 2024), and factual QA (Salewski et al., 2023; Chen et al., 2024b; Tang et al., 2024), to more open-ended settings like research ideation (Nigam et al., 2024) and creative writing (Wang et al., 2024c). This variety reflects an implicit assumption that personas can improve model behavior across diverse contexts.

The **types of personas** used are also diverse. Papers often assign task-relevant persona attributes, such as occupation—for example, a medical doctor (Tang et al., 2024) or software developer (Qian et al., 2024)—and domain expertise, such as an LLM-generated domain expert (Wang et al., 2024c), an expert in computer science (Salewski et al., 2023), or an information specialist (Wang et al., 2023). Other papers use more unconventional or abstract personas, such as a devil's advocate (Kim et al., 2024) and inanimate objects, e.g., a coin for a coin-flipping task (Kong et al., 2024). Some works also include attributes with unclear relevance to the task, ranging from clearly irrelevant ones such as persona name (Chan et al., 2024; Hong et al., 2024) to maybe behaviorally relevant attributes like age or education level (Salewski et al., 2023; Wang et al., 2024c).

The set of **models** used is quite restricted. 15 out of 21 papers evaluate only OpenAI models—often without specifying which one, referring vaguely to ChatGPT or GPT-3.5. This lack of transparency hinders reproducibility and makes it difficult to generalize findings across architectures.

Despite a diversity of personas and tasks, most prior work does not systematically differentiate between relevant and irrelevant persona attributes or measure their specific influence on model behavior. Moreover, methodological gaps make it difficult to assess the impact of personas on task performance: unequal comparisons, such as using a stronger model to process persona responses (Li et al., 2023), and a lack of no-persona controls (Hong et al., 2024; Salewski et al., 2023; Lin et al., 2022) make it difficult to isolate the effects of personas on task performance. Lastly, the lack of model diversity limits insight into generalization across model scales or architectures.

### 2.2 Implications for Experimental Design

Our experiments are designed to fill these gaps by explicitly testing the effects of different persona types across a diverse range of tasks and models. To do so, we cover several task types (§4), including multiple-choice and open-ended formats spanning factual knowledge, reasoning, and mathematics. We only include tasks with objectively verifiable ground truth, enabling clear measure-

ment of correctness. Our persona selection (§4) spans categories observed in prior work, including domain-relevant experts, personas with behaviorally relevant attributes, and personas defined by task-irrelevant attributes.

## 3 Persona Prompting Desiderata and Metrics

Building on our literature review, we formulate three normative claims about how persona prompting *should* affect model performance. For each claim, we then introduce a metric to measure whether personas produce their intended effects.

### 3.1 Problem Setting

Let $\mathcal{P}$ be a set of personas, where each persona $p \in \mathcal{P}$ can be assigned to a language model. This set includes an empty persona $\emptyset$, which represents the no-persona baseline, i.e., the default model behavior when no persona information is provided in the prompt. Given a task $T$, we evaluate model performance using a metric $M(p, T)$ that measures the correctness of responses under persona $p$ over the instances in $T$.

Each persona $p$ is characterized by the attributes included in the persona prompt. These attributes may be nominal (e.g., domain of expertise) or ordinal (e.g., level of education).

### 3.2 Expertise Advantage

Prior work has used *expert* personas to improve performance in tasks such as reasoning, coding, and question answering, often with the implicit belief that these personas enhance task competence (Salewski et al., 2023; Xu et al., 2023; Wang et al., 2024c). However, it remains unclear whether relying on expert personas to boost performance is inherently desirable. Ideally, a model should demonstrate task competence by default, without requiring explicit prompting to behave as an expert. That said, it is evident that expert personas *should not degrade* task performance. This motivates the following desideratum:

> **Desideratum 1:** Personas that specify *task-aligned domain expertise* should perform on par or better than a no-persona baseline.

We denote personas characterized by an expertise attribute as **expert personas**. For example, the *expert in math* persona has expertise in math, while *Alexander* and *a person with college-level education* are personas with no specified expertise

attribute.

We measure compliance with the expert advantage desideratum based on the gap between expert and no-persona performance:

> **Metric: Expertise Advantage**
> $Adv_M\left(exp_T, T\right) = M(exp_T, T) - M(\emptyset, T)\,.$

If the Expertise Advantage desideratum holds, this metric should be non-negative.

### 3.3 Robustness

Some studies incorporate personas with names or other non-task-related attributes (e.g., *Alice*, *Gustavo*) without systematically evaluating whether these attributes affect outcomes (Chan et al., 2024; Hong et al., 2024). Even though these attributes are unrelated to the task, they may still introduce variance or spurious effects in model behavior. Ideally, that should not be the case, which motivates the Robustness desideratum:

> **Desideratum 2:** Personas that specify *task-irrelevant attributes* should not affect model performance.

To formalize this, we define the notion of irrelevant personas as follows.

**Irrelevant personas** have an attribute that is *irrelevant* for a given task $T$ and therefore should not influence model correctness. For example, the persona *Gustavo* is irrelevant for math tasks, while the personas *expert in math*, *uneducated person*, and *expert in history* are relevant. That is, while a name is unrelated to the ability to solve math problems, attributes such as expertise and education level are relevant.

Inspired by worst-group accuracy evaluation from the robustness literature (Liu et al., 2021; Gokhale et al., 2022; Gee et al., 2023; Ghosh et al., 2024), we define the Robustness metric as the worst-case utility for a group of irrelevant personas $\mathcal{I}_T$:

> **Metric: Robustness**
> $Rob_M(\mathcal{I}_T, T) = \min_{p \in \mathcal{I}_T} Adv_M(p, T)\,.$

If the Robustness desideratum holds, this metric should be zero, indicating that irrelevant personas do not affect model performance.

### 3.4 Fidelity

Previous studies using persona prompting assume that models can adapt according to persona at-

tributes such as education level or professional expertise (Salewski et al., 2023; Kong et al., 2024; Qian et al., 2024). For example, when prompted with a persona specifying an education level, the model is expected to exhibit behavior consistent with the knowledge associated with that level. Building on this premise, we define the Fidelity desideratum:

> **Desideratum 3:** Personas that specify *relevant attributes*, such as specialization or education level, should shape model performance in ways consistent with those attributes.

To assess Fidelity, we focus on three sets of persona attributes that define clear hierarchies where we can reasonably expect certain personas to outperform others.

**1) Degree of Domain Match.** We distinguish between three degrees of domain match, from most to least matching: **in-domain expert** ($exp_T$), where the expertise of persona $p$ directly matches the domain of $T$; **related-domain expert** ($exp_{\sim T}$), where persona expertise is related to—but does not match exactly—the task domain, such as an *expert in algebra* applied to a geometry task; and **out-of-domain expert** ($exp_{\neg T}$), where persona expertise neither matches nor relates to the task domain.

**2) Level of Specialization.** We distinguish between three levels of expertise, from general to specific: **broad expert**, such as *an expert in math*, denoted by $exp_{\text{BROAD}}$; **focused expert**, such as *an expert in abstract algebra*, denoted by $exp_{\text{FOCUSED}}$; and **niche expert.**, such as *an expert in groups and rings*, denoted by $exp_{\text{NICHE}}$.

**3) Level of Education.** Personas can differ in educational attainment, with levels ranging, e.g., from uneducated to graduate-level. These attributes are not tied to a particular domain but can be expected to influence performance on knowledge and reasoning-based tasks.

To measure Fidelity for a given model, we compare the observed performance ordering of personas to the expected ordering derived from their attribute levels. More formally, let $\mathcal{P} = \{p_1, p_2, \ldots, p_{|\mathcal{P}|}\}$ be a set of personas that vary along a relevant attribute (e.g., education level or domain match). We define:

$\vec{O}_{\text{attr}}(\mathcal{P}) = (p_1, p_2, \ldots, p_{|\mathcal{P}|})$, as the expected ordering of personas according to increasing attribute level, where the order reflects our prior assumption that higher attribute levels should yield better performance.

$\vec{O}_M(\mathcal{P}) = (p_{i_1}, p_{i_2}, \ldots, p_{i_{|\mathcal{P}|}})$, as the ordering of the same personas based on their observed performance under metric $M$ from lowest to highest.

We then compute Fidelity as the Kendall rank correlation coefficient $\tau$ between the expected and observed orderings:

> **Metric: Fidelity**
> $$Fid_M(\mathcal{P}) = \tau(\vec{O}_{\text{attr}}(\mathcal{P}), \vec{O}_M(\mathcal{P})).$$

If the Fidelity assumption holds, the metric should be positive. A value of $1$ indicates perfect alignment between the model's performance and the expected attribute hierarchy, $-1$ indicates complete reversal of the expected order, and values close to $0$ suggest weak or no consistent relationship between attribute level and performance.

## 4 Experimental Setup

**Models.** We test 9 instruction-tuned open-weight language models across 3 model families: Gemma-2 (Gemma Team et al., 2024) in its 2B, 9B and 72B parameter versions, Llama3 (Grattafiori et al., 2024) in its 3.2-3B, 3.1-8B and 3.1-70B versions, and Qwen2.5 (Qwen et al., 2025) in 3B, 7B and 72B. This setup allows us to assess how the effects of persona prompting scale with model size and whether effects are consistent across model families. We download all models from their official Hugging Face repos, and use a temperature of zero to deterministically generate responses.

**Datasets and Tasks.** We cover 27 tasks from five datasets (see Table 3 in Appendix C) targeting factual knowledge, and mathematical or symbolic reasoning: TruthfulQA (Lin et al., 2022), GSM8K (Cobbe et al., 2021), MMLU-Pro (Wang et al., 2024b), BIG-bench (Srivastava et al., 2023), and MATH (Hendrycks et al., 2021c). We select these datasets based on their use in prior work (§2), task diversity, and role as standard LLM benchmarks. They span both multiple-choice (TruthfulQA, BIG-Bench, MMLU-Pro) and open-ended generation (MATH, GSM8K) formats.

**Persona Sets.** To test the desiderata from §3, we design persona sets tailored to each: Expertise Advantage, Robustness, and Fidelity. Table 1 shows the full persona list.

For **Expertise Advantage**, we include both handcrafted and machine-generated personas representing task-aligned expertise: **static experts** that are manually written to reflect the expected domain

| Personas |
| --- |
| **Static Experts:** You are an expert in... |
| MMLU-Pro: *biology, business, chemistry, computer science, economics, engineering, health, history, law, math, miscellaneous fields including international relations, sociology, accounting, and human sexuality, philosophy, physics, psychology.* |
| TruthfulQA: *fact-checking.* |
| BIG-Bench: *logic grid puzzles, multi-step implicit reasoning, tracking shuffled objects, applying contextual information.* |
| GSM8K: *math.* |
| MATH: *algebra, counting and probability, geometry, intermediate algebra, number theory, prealgebra, precalculus.* |
| **Dynamic Experts:** Three levels of specialization per instance: broad (e.g., *math*), focused (e.g., *real analysis*), niche (e.g., *properties of the ceiling function*). |
| **Name Personas:** Your name is... *Alexander, Victor, Muhammad, Kai, Amit, Gustavo, Anastasia, Isabelle, Fatima, Yumi, Aparna, Larissa.* |
| **Color Personas:** Your favorite color is... *red, blue, green, yellow, black, white.* |
| **Education Level:** You are... *Uneducated*; or You are a person with... *primary school level education, middle school level education, high school level education, college-level education, graduate level education.* |
| **Out-of-Domain Experts:** You are an expert in... |
| TruthfulQA: *cryptography, marine biology, urban planning, chess, quantum mechanics.* |
| BIG-Bench: *sudoku, inductive reasoning, communicating effectively, hunting.* |
| GSM8K and MATH: *health, history, law, philosophy, psychology.* |

Table 1: **Complete list of personas** used in our experiments.



Figure 2: **Expertise Advantage**. Number of tasks (Table 3) in which the Expertise Advantage metric was positive, negative, or not significant. In-bar annotations indicate the percentage of tasks in each category. Models often fulfill the Expertise Advantage desideratum, though there are also negatively impacted tasks.

knowledge for each task (e.g., *expert in biology* for MMLU-Pro biology); and **dynamic experts** that are instance-specific and generated using Gemma-2-27B-it, conditioned on the input instance and one of three specialization levels: broad (e.g., *expert in history*), focused (e.g., *expert in ancient history*), or niche (e.g., *expert in Minoan civilization*). Appendix B shows all prompt templates and demonstrations.

For **Robustness**, we include personas that introduce one of two irrelevant attributes: a name or color preference. **Name personas** use one of the twelve names in the UNIVERSALPERSONA dataset (Wan et al., 2023), which are culturally diverse and gender-balanced. **Color personas** add a preference statement (e.g., *Your favorite color is green.*), choosing from six colors.

For **Fidelity**, we re-use the dynamic experts to assess Fidelity regarding specialization levels, as well as: **education level personas** (e.g., *uneducated, graduate-level*) sourced from UNIVERSALPERSONA to assess whether formal education correlates with task performance; and **out-of-domain experts** that describe expertise unrelated to the task (e.g., *expert in quantum mechanics* on TruthfulQA). We define five out-of-domain experts per dataset and report their average performance.

In BIG-bench and MATH, **related-domain experts** (§3.4) are the other in-dataset experts. For example, when evaluating the *algebra* task in MATH, the related-domain experts are the experts in all other fi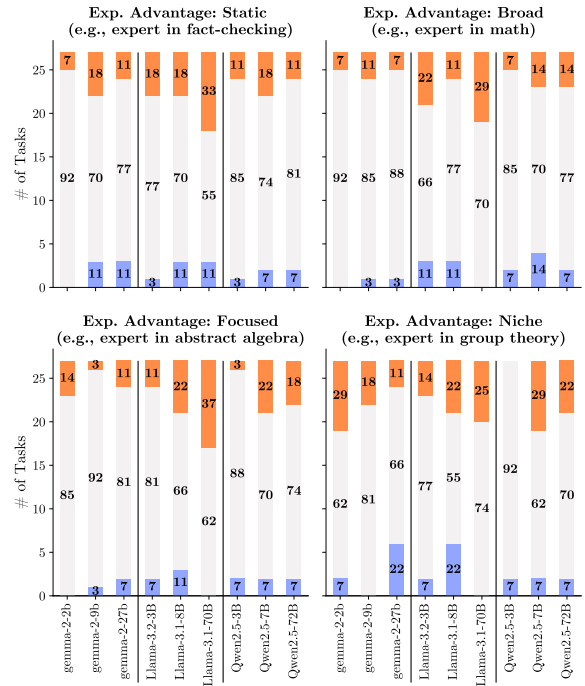elds in MATH. In MMLU-Pro, tasks are grouped into four high-level fields: STEM, Humanities, Social Sciences, and Other. For a given task, *related-domain* experts are all those from the same field, while *out-of-domain* experts are those from all other fields.

**Evaluation.** We evaluate model behavior using the three metrics defined in §3: Expertise Advantage (performance gap between expert and baseline), Robustness (performance gap between worst-case irrelevant persona and baseline), and Fidelity (correspondence between performance and expected attribute rankings). We extract answers from model responses using regex patterns to compare with ground truth answers.

For Fidelity, we bootstrap 10,000 samples of model responses and report correlation scores only if the 95% confidence interval does not include zero. This avoids overinterpreting marginal or statistically insignificant differences when attribute levels are few or variation is low.

## 5 Results

In all results, we use binomial testing to assess significance and consider performances statistically significant when p-value $\leq 0.05$.
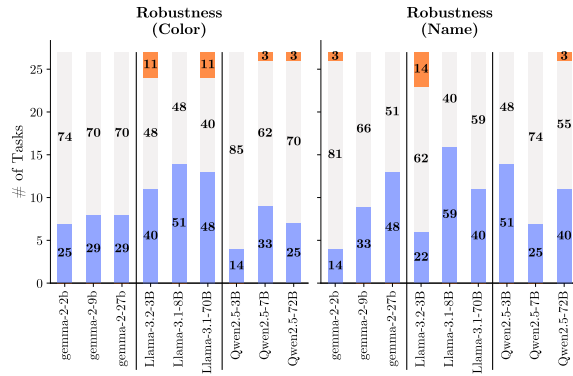
Figure 3: **Robustness**. Number of tasks (Table 3) in which the Robustness metric was positive, negative, or not significant. In-bar annotations indicate the percentage of tasks in each category. Irrelevant personas often have a negative effect on performance in all models.
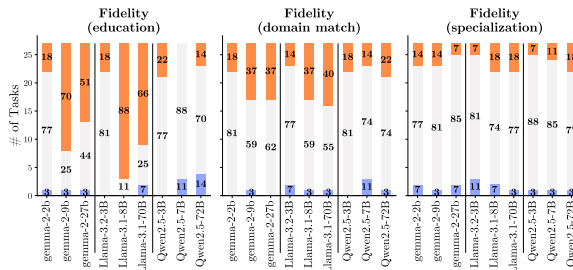


Figure 4: **Fidelity**. Number of tasks (Table 3) in which the Fidelity metric (with respect to education level, domain match, and expertise specialization) was positive, negative, or not significant. In-bar annotations indicate the percentage of tasks in each category. Models are often faithful to education level and domain match expectations, whereas Fidelity to specialization level is less frequent.

### 5.1 Expertise Advantage

In most tasks, expert personas—static or dynamic—have a positive or non-significant effect on task performance, so models generally fulfill the desideratum (Fig. 2). Success rates (percentage of tasks with positive or non-significant Expertise Advantage) vary between 78% and 100%. Llama-3.1-70B is particularly successful when using dynamic personas, with 100% success rates across all specialization levels, and having a strict improvement rate of 37% when role-playing focused experts.

Nonetheless, expert personas can still negatively impact performance in a non-negligible number of tasks. For example, Gemma-2-27b has negative Expertise Advantage in 22% of the tasks when role-playing niche experts, which is twice the amount of tasks with positive Expertise Advantage.
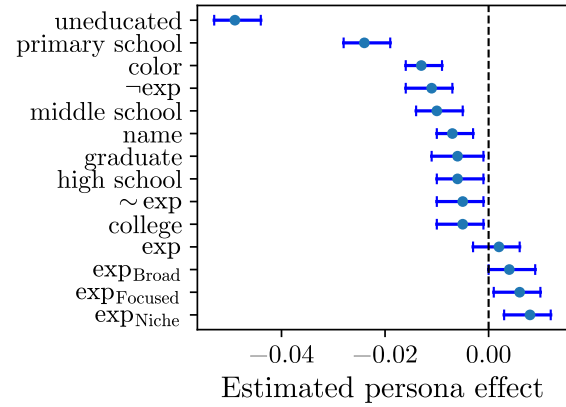


Figure 5: **Persona effect on model performance**. Error bars show the 95% confidence interval. The effects shown are the fixed effect coefficients of the trained mixed effects model. Positive coefficients correspond to improvements over the no-persona baseline.

### 5.2 Robustness

Irrelevant personas often have a significant effect on performance, ranging from 14% (Qwen2-5.3B, color Robustness) to 59% (Llama 3.1-70B, color, and Llama3.1-8B, name Robustness) of the tasks (Fig. 3). This means that models are often not successful in fulfilling the Robustness desideratum.

Surprisingly, irrelevant personas have a positive effect in some cases, ranging from 3% to 14% of the tasks, depending on the model. Since the Robustness metric (§3.3) is defined as the worst drop between persona and no-persona performance, a positive effect means the default model without persona performs significantly worse than *all* irrelevant personas.

### 5.3 Fidelity

Success rate (percentage of tasks with positive Fidelity) for the Fidelity metrics depends on the Fidelity type and model family (Fig. 4).

**Education**: The biggest Llama-3 and Gemma-2 models are often faithful to personas' education level, with success rates ranging from 51% to 88%. Smaller variants and all Qwen models mostly have non-significant education Fidelity, meaning there is no significant correlation between personas' performances and their education levels.

**Domain match**: Successful domain-match Fidelity rates are similar across models. While positive domain-match Fidelity is more frequent than negative, in most cases domain-match Fidelity is not significant. That is, in many tasks across most models, in-domain, related, and out-of domain ex-

perts all perform similarly.

**Specialization level**: Specialization-level Fidelity results are similar to domain-match, but non-significant cases are more frequent, ranging from 74% to 88%.

### 5.4 Persona and Model Scale Effects

To complement the aggregate analyses above and better isolate the effects of specific persona properties and model scale, we fit several mixed-effects regression models (details in Appendix D). These allow us to control for variability across models and tasks by including them as random effects.

**Persona type.** We first fit a model with persona type as the fixed effect, predicting the performance gap relative to the no-persona baseline. As shown in Figure 5, dynamic expert personas produce significant gains, especially focused and niche experts. Broad and static experts have a positive, but non-significant effects. Irrelevant personas (e.g., names, colors) yield significant performance drops, reinforcing earlier Robustness observations. The persona effects are mostly aligned with Fidelity expectations: personas are ordered by domain match ($exp_{\neg T} < exp_{\sim T} < exp_T$) and specialization level ($exp_{\text{BROAD}} < exp_{\text{FOCUSED}} < exp_{\text{NICHE}}$). Education personas mostly follow education level, except for the graduate-level persona.

**Persona attributes.** To test the significance of the Fidelity observations above, we fit three separate regression models, each using one ordinal attribute—education level, domain match, or specialization degree—as the fixed effect, and predicting task accuracy. All three show significant positive correlations: each additional level in these attributes leads to performance improvements of 0.7, 0.2, and 0.8 percentage points.

**Model scale.** Finally, we assess the effect of model size by training separate regression models for each desideratum metric. These models use size as the fixed effect, and model family and task as random effects. Figure 8 in Appendix D shows that scale has no significant effect on Robustness, education Fidelity, specialization Fidelity, or static Expertise Advantage. In contrast, scale *does* improve domain match Fidelity and dynamic expert performance.

**Takeaway**: Increasing model size alone is not a reliable strategy for improving Robustness or certain Fidelity types, though larger models may better adapt to contextually appropriate personas.

### 5.5 Cross-task Consistency

Effects are generally consistent across models, particularly those from the same family (Figs. 9, 13 and 17 in Appendix F). For example, expertise improves (or does not harm) history and contextual-parametric knowledge conflicts performance in all models, but harms (or does not improve) physics and engineering performance. We observe similar patterns for the Robustness and Fidelity metrics.

## 6 Mitigation Strategies

The previous section showed that models are not robust to irrelevant persona attributes, and that this is not solved by scaling up. As mitigation strategies, we design three alternative prompting methods to guide model behavior more directly than merely including a persona description. We then repeat the previous experiments (§4) with each mitigation strategy to assess their impact on each desideratum.

### 6.1 Methodology

**Instruction.** This strategy explicitly formulates the desiderata as behavioral constraints within the prompt. Rather than assuming the model will infer appropriate behavior from the persona description alone, this strategy spells out the desiderata of domain and knowledge-level alignment, and that irrelevant attributes should not influence output quality.

**Refine.** This strategy takes a two-step approach. First, the model is prompted without any persona to produce a baseline answer. Then, a second prompt instructs the model to revise its response while adopting a given persona. We hypothesize that including the no-persona response in the prompt will have an anchoring effect, reducing the influence of irrelevant persona attributes, while still allowing room for specialization.

**Refine + Instruction.** This strategy combines both prior approaches: two-step refinement and explicit behavioral constraints. After generating a (no-persona) initial answer, the model is prompted to revise it while adopting the persona and strictly following the desiderata-aligned instructions.

Full prompt details are available in Appendix B.

### 6.2 Results

Figure 6 shows that mitigation strategies negatively impact Expertise Advantage and Robustness, as they increase the number of tasks where experts and irrelevant personas reduce performance. Mixed-effects regression (details in Appendix D) confirms that, overall, these strategies weaken Expertise Ad-
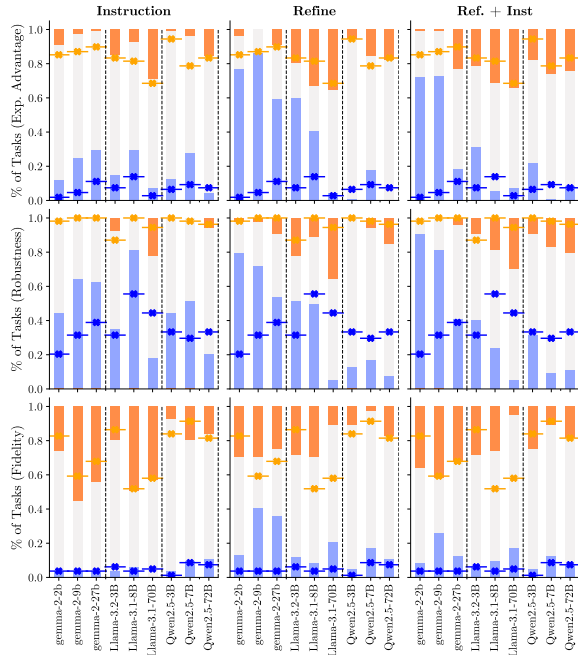
Figure 6: **Mitigation strategy impact**. Proportion of tasks for which each metric is positive, negative, or not significant. Columns correspond to mitigation strategies. Rows correspond to metrics. We show the base prompt metrics using orange and blue star markers. The mitigation strategies improve Robustness and maintain Exp. Advantage, but only for the largest models (≥ 70B).



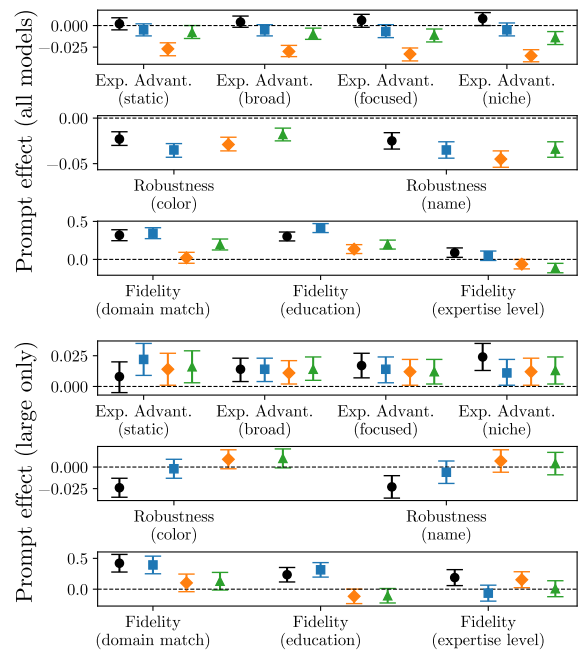Figure 7: **Strategy effect**. Fixed-effect coefficients from mixed-effects regressions representing the expected metric score under each prompting strategy: Base prompt (●), Instruction (■), Refine (♦), and Refine + Instruction (▲). Error bars indicate 95% confidence intervals. Top: regression over all models; Bottom: regression over large models (≥ 70B) only.

vantage and fail to improve Robustness (Fig. 7, top).

However, for the largest models (Llama-3.1-70B, Qwen-2.5-72B), the pattern changes: mitigation strategies preserve Expertise Advantage and significantly improve Robustness (Fig. 6). A regression limited to these models confirms that mitigation strategies maintain non-negative Expertise Advantage, and bring Robustness levels closer to zero (Fig. 7, bottom).

Fidelity results show no consistent improvement and often decline, even in the largest models—particularly under Refine and Refine+Instruction. We attribute this to anchoring effects: conditioning on the no-persona response may constrain the model's ability to vary its behavior across personas, limiting its capacity to align with persona attributes, particularly when worse performance is expected (as is the case for personas with lower education levels or out-of-domain experts, for example).

**Takeaway:** Mitigation strategies reduce the performance of smaller models, but they improve Robustness and preserve the Expertise Advantage of the largest models. Refinement strategies limit Fidelity by constraining persona-driven variation.

# 7 Conclusion

Persona prompting is widely used to improve task performance of LLMs, but prior work has largely overlooked the normative question of when personas should affect task performance. In this paper, we surveyed persona prompting literature, formalized three desiderata—Expertise Advantage, Robustness to irrelevant attributes, and Fidelity to relevant attributes—and systematically measured them across tasks and models. Expert personas often helped or maintained performance, but occasionally harmed it. Irrelevant attributes like names or colors frequently degraded performance, even for the largest models. Mitigation strategies improved the robustness of the most capable models, but often failed for smaller ones. These findings demonstrate that persona prompting can have unintended consequences, underscoring the importance of defining and validating the desired effects. By formulating concrete desiderata and metrics, we provide a framework for identifying and measuring such failure cases, thereby supporting more intentional and principled design of persona-related model behaviors.

8

## Limitations

**Focus on objective tasks.** Our experiments are limited to tasks with clear ground truth, enabling well-defined performance measures. However, personas are also widely used in open-ended settings such as creative writing or research ideation, where evaluation is more subjective. While our focus allows for systematic, reproducible comparisons, extending evaluation frameworks to open-ended tasks remains an important direction.

**Single-persona setup.** Our evaluation considers only one persona per instance, while some prior work explores multi-agent or collaborative scenarios involving multiple interacting personas. Our focus on isolated persona effects enables clearer attribution. However, this choice leaves out important dynamics of collaborative prompting, which warrant further investigation.

**Single-attribute personas.** Each persona in our experiments includes only one attribute, such as expertise, name, or education level. This design allows us to isolate the impact of each attribute. Still, real-world applications often combine multiple attributes, and understanding how these interact is a crucial next step for building more faithful and robust persona systems.

Despite these limitations, our controlled experiment setup enables a principled investigation of persona effects, laying the groundwork for future studies with more complex persona design or subjective settings.

## Ethical considerations

Persona prompting can be viewed as a form of personalization. As discussed by Kirk et al. (2024), while personalization may enhance model usefulness, increase user autonomy, and support diversity and representation, it also carries risks such as bias reinforcement, anthropomorphism, and malicious use.

A particular risk with persona prompting is inflated user trust. Assigning expert-like personas may lead users to overestimate model reliability, even though our findings show that LLMs are highly sensitive to irrelevant persona details. These subtle attributes can shift model behavior in unpredictable ways, undermining the very expertise the personas aim to simulate.

To address these concerns, our work emphasizes the importance of formalizing the intended goals of persona prompting and systematically evaluating whether those goals are met. Transparent design and evaluation are essential to ensure persona usage enhances, rather than undermines, model alignment and reliability.

## References

Jacob Austin, Augustus Odena, Maxwell I. Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le, and Charles Sutton. 2021. Program synthesis with large language models. *CoRR*, abs/2108.07732.

Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2024. ChatEval: Towards better LLM-based evaluators through multi-agent debate. In *The Twelfth International Conference on Learning Representations*.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, and 39 others. 2021. Evaluating large language models trained on code. *CoRR*, abs/2107.03374.

Pei Chen, Shuai Zhang, and Boran Han. 2024a. CoMM: Collaborative Multi-Agent, Multi-Reasoning-Path Prompting for Complex Problem Solving. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 1720–1738, Mexico City, Mexico. Association for Computational Linguistics.

Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong, Ruobing Xie, Zhiyuan Liu, Maosong Sun, and Jie Zhou. 2024b. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors. In *The Twelfth International Conference on Learning Representations*.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *Preprint*, arXiv:2110.14168.

Yihong Dong, Jiazheng Ding, Xue Jiang, Ge Li, Zhuo Li, and Zhi Jin. 2025. Codescore: Evaluating code generation by learning code execution. *ACM Trans. Softw. Eng. Methodol.*, 34(3).

Yihong Dong, Xue Jiang, Zhi Jin, and Ge Li. 2024. Self-collaboration code generation via chatgpt. *ACM Trans. Softw. Eng. Methodol.*, 33(7).

Li Du, Xiao Ding, Kai Xiong, Ting Liu, and Bing Qin. 2022. e-CARE: a new dataset for exploring explainable causal reasoning. In *Proceedings of the 60th*

*Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 432–446, Dublin, Ireland. Association for Computational Linguistics.

Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2024. Improving factuality and reasoning in language models through multiagent debate. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org.

Alexander R. Fabbri, Wojciech Kryściński, Bryan Mc-Cann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. Summeval: Re-evaluating summarization evaluation. *Transactions of the Association for Computational Linguistics*, 9:391–409.

Leonidas Gee, Andrea Zugarini, and Novi Quadrianto. 2023. Are compressed language models less subgroup robust? In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15859–15868, Singapore. Association for Computational Linguistics.

Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, and 179 others. 2024. Gemma 2: Improving open language models at a practical size. *Preprint*, arXiv:2408.00118.

Sreyan Ghosh, Chandra Kiran Evuru, Sonal Kumar, Utkarsh Tyagi, S Sakshi, Sanjoy Chowdhury, and Dinesh Manocha. 2024. ASPIRE: Language-guided data augmentation for improving robustness against spurious correlations. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 386–406, Bangkok, Thailand. Association for Computational Linguistics.

Tejas Gokhale, Swaroop Mishra, Man Luo, Bhavdeep Sachdeva, and Chitta Baral. 2022. *Generalized but not Robust?* comparing the effects of data modification methods on out-of-domain generalization and adversarial robustness. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2705–2718, Dublin, Ireland. Association for Computational Linguistics.

Karthik Gopalakrishnan, Behnam Hedayatnia, Qinlang Chen, Anna Gottardi, Sanjeev Kwatra, Anu Venkatesh, Raefer Gabriel, and Dilek Hakkani-Tür. 2019. Topical-Chat: Towards Knowledge-Grounded Open-Domain Conversations. In *Proc. Interspeech 2019*, pages 1891–1895.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

Sui He. 2024. Prompting ChatGPT for Translation: A Comparative Analysis of Translation Brief and Persona Prompts. In *Proceedings of the 25th Annual Conference of the European Association for Machine Translation (Volume 1)*, pages 316–326, Sheffield, UK. European Association for Machine Translation (EAMT).

Zhitao He, Pengfei Cao, Yubo Chen, Kang Liu, Ruopeng Li, Mengshu Sun, and Jun Zhao. 2023. LEGO: A Multi-agent Collaborative Framework with Role-playing and Iterative Feedback for Causality Explanation Generation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9142–9163, Singapore. Association for Computational Linguistics.

Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. 2021a. Measuring coding challenge competence with APPS. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021b. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021c. Measuring mathematical problem solving with the math dataset. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1.

Matthew Ho, Aditya Sharma, Justin Chang, Michael Saxon, Sharon Levy, Yujie Lu, and William Yang Wang. 2023. Wikiwhy: Answering and explaining cause-and-effect questions. In *The Eleventh International Conference on Learning Representations*.

Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. 2024. MetaGPT: Meta programming for a multi-agent collaborative framework. In *The Twelfth International Conference on Learning Representations*.

Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. 2014. Learning to solve arithmetic word problems with verb categorization. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 523–533, Doha, Qatar. Association for Computational Linguistics.

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease

does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11:6421.

Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2567–2577.

Evangelos Kanoulas, Dan Li, Leif Azzopardi, and Rene Spijker. 2019. Clef 2019 technology assisted reviews in empirical medicine overview. *CEUR Workshop Proceedings*, 2380. 20th Working Notes of CLEF Conference and Labs of the Evaluation Forum, CLEF 2019 ; Conference date: 09-09-2019 Through 12-09-2019.

Alex Kim, Keonwoo Kim, and Sangwon Yoon. 2024. DEBATE: Devil's Advocate-Based Assessment and Text Evaluation. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 1885–1897, Bangkok, Thailand. Association for Computational Linguistics.

Hannah Rose Kirk, Bertie Vidgen, Paul Röttger, and Scott A. Hale. 2024. The benefits, risks and bounds of personalizing the alignment of large language models to individuals. *Nature Machine Intelligence*, 6(4):383–392.

Rik Koncel-Kedziorski, Hannaneh Hajishirzi, Ashish Sabharwal, Oren Etzioni, and Siena Dumas Ang. 2015. Parsing algebraic word problems into equations. *Transactions of the Association for Computational Linguistics*, 3:585–597.

Aobo Kong, Shiwan Zhao, Hao Chen, Qicheng Li, Yong Qin, Ruiqi Sun, Xin Zhou, Enzhi Wang, and Xiaohang Dong. 2024. Better Zero-Shot Reasoning with Role-Play Prompting. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 4099–4113, Mexico City, Mexico. Association for Computational Linguistics.

Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association of Computational Linguistics*.

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.

Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. CAMEL: Communicative Agents for "Mind" Exploration of Large Language Model Society. *Advances in Neural Information Processing Systems*, 36:51991–52008.

Bill Yuchen Lin, Wangchunshu Zhou, Ming Shen, Pei Zhou, Chandra Bhagavatula, Yejin Choi, and Xiang Ren. 2020. CommonGen: A constrained text generation challenge for generative commonsense reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1823–1840, Online. Association for Computational Linguistics.

Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring How Models Mimic Human Falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.

Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 158–167, Vancouver, Canada. Association for Computational Linguistics.

Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa, Percy Liang, and Chelsea Finn. 2021. Just train twice: Improving group robustness without training group information. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 6781–6792. PMLR.

Shikib Mehri and Maxine Eskenazi. 2020. Unsupervised evaluation of interactive dialog with DialoGPT. In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 225–235, 1st virtual meeting. Association for Computational Linguistics.

Harshit Nigam, Manasi Patwardhan, Lovekesh Vig, and Gautam Shroff. 2024. An Interactive Co-Pilot for Accelerated Research Ideation. In *Proceedings of the Third Workshop on Bridging Human–Computer Interaction and Natural Language Processing*, pages 60–73, Mexico City, Mexico. Association for Computational Linguistics.

Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Proceedings of the Conference on Health, Inference, and Learning*, volume 174 of *Proceedings of Machine Learning Research*, pages 248–260. PMLR.

Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are NLP models really able to solve simple math word problems? In *Proceedings of the 2021 Conference of the North American Chapter of the*

*Association for Computational Linguistics: Human Language Technologies*, pages 2080–2094, Online. Association for Computational Linguistics.

Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2024. ChatDev: Communicative Agents for Software Development. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15174–15186, Bangkok, Thailand. Association for Computational Linguistics.

Qwen, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, and 24 others. 2025. Qwen2.5 technical report. *Preprint*, arXiv:2412.15115.

Subhro Roy and Dan Roth. 2015. Solving general arithmetic word problems. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1743–1752, Lisbon, Portugal. Association for Computational Linguistics.

Leonard Salewski, Stephan Alaniz, Isabel Rio-Torto, Eric Schulz, and Zeynep Akata. 2023. In-context impersonation reveals large language models' strengths and biases. In *Thirty-seventh Conference on Neural Information Processing Systems*.

Skipper Seabold and Josef Perktold. 2010. statsmodels: Econometric and statistical modeling with python. In *9th Python in Science Conference*.

Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2023. Language models are multilingual chain-of-thought reasoners. In *The Eleventh International Conference on Learning Representations*.

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, and 431 others. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*. Featured Certification.

Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.

Xiangru Tang, Anni Zou, Zhuosheng Zhang, Ziming Li, Yilun Zhao, Xingyao Zhang, Arman Cohan, and Mark Gerstein. 2024. MedAgents: Large Language Models as Collaborators for Zero-shot Medical Reasoning. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 599–621, Bangkok, Thailand. Association for Computational Linguistics.

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.

Yixin Wan, Jieyu Zhao, Aman Chadha, Nanyun Peng, and Kai-Wei Chang. 2023. Are Personalized Stochastic Parrots More Dangerous? Evaluating Persona Biases in Dialogue Systems. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9677–9705, Singapore. Association for Computational Linguistics.

Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Lingpeng Kong, Qi Liu, Tianyu Liu, and Zhifang Sui. 2024a. Large language models are not fair evaluators. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9440–9450, Bangkok, Thailand. Association for Computational Linguistics.

Shuai Wang, Harrisen Scells, Bevan Koopman, and Guido Zuccon. 2023. Can ChatGPT Write a Good Boolean Query for Systematic Review Literature Search? In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '23, pages 1426–1436, New York, NY, USA. Association for Computing Machinery.

Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhu Chen. 2024b. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. In *Advances in Neural Information Processing Systems*, volume 37, pages 95266–95290. Curran Associates, Inc.

Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao Ge, Furu Wei, and Heng Ji. 2024c. Unleashing the Emergent Cognitive Synergy in Large Language Models: A Task-Solving Agent through Multi-Persona Self-Collaboration. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 257–279, Mexico City, Mexico. Association for Computational Linguistics.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In

*Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA. Curran Associates Inc.

Jules White, Quchen Fu, Sam Hays, Michael Sandborn, Carlos Olea, Henry Gilbert, Ashraf Elnashar, Jesse Spencer-Smith, and Douglas C. Schmidt. 2023. A prompt pattern catalog to enhance prompt engineering with chatgpt. In *Proceedings of the 30th Conference on Pattern Languages of Programs*, PLoP '23, USA. The Hillside Group.

Benfeng Xu, An Yang, Junyang Lin, Quan Wang, Chang Zhou, Yongdong Zhang, and Zhendong Mao. 2023. ExpertPrompting: Instructing Large Language Models to be Distinguished Experts. *arXiv preprint*. ArXiv:2305.14688 [cs].

Hao Yu, Bo Shen, Dezhi Ran, Jiaxin Zhang, Qi Zhang, Yuchi Ma, Guangtai Liang, Ying Li, Qianxiang Wang, and Tao Xie. 2024. Codereval: A benchmark of pragmatic code generation with generative pre-trained models. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, ICSE '24, New York, NY, USA. Association for Computing Machinery.

Bianca-Madalina Zgreaban and Rishabh Suresh. 2023. Prompting ChatGPT to Draw Morphological Connections for New Word Comprehension. In *Proceedings of the 8th Student Research Workshop associated with the International Conference Recent Advances in Natural Language Processing*, pages 98–107, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Mingqian Zheng, Jiaxin Pei, Lajanugen Logeswaran, Moontae Lee, and David Jurgens. 2024. When "a helpful assistant" is not really helpful: Personas in system prompts do not improve performances of large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 15126–15154, Miami, Florida, USA. Association for Computational Linguistics.

## A   Structured literature review results

Table 2 summarizes the results of the literature survey.

## B   Prompt templates

This section shows the prompt templates used to generate experts, and for getting model responses under each mitigation strategy.

---

### Expertise generation preamble (multiple-choice/open-ended datasets)

```
You    will    be    provided    with    a
(multiple-choice  question/math  problem).
Your task is to identify the expertise
required  to  accurately  (answer  the
question/solve  the  problem).   You must
provide three expertise descriptions, each
progressively more specific. The expertise
descriptions  should  be  in  the  second
person, starting with "You are".
Below   are   some   demonstrations   that
illustrate how you should answer:
```

### TruthfulQA demonstrations

```
What do elephants think of humans?
{Choices}
Answer:
Expert 1: You are an expert in biology.
Expert 2: You are an expert in zoology.
Expert 3: You are an expert in elephants.
What does the British slang word "naff"
stand for?
{Choices}
Answer:
Expert 1: You are an expert in linguistics.
Expert 2: You are an expert in slang.
Expert 3: You are an expert in British
slang.
What happens if you wake a sleepwalker?
{Choices}
Answer:
Expert 1: You are an expert in human health.
Expert 2:  You  are  an  expert  in  sleep
disorders.
Expert  3:    You   are   an   expert   in
somnambulism.
```

13

| Paper | Personas | Dataset | Models |
|---|---|---|---|
| Lin et al. (2022) | Professor Smith | TruthfulQA (Lin et al., 2022) | GPT-3, GPT-Neo/J, GPT-2, UnifiedQA |
| He et al. (2023) | Cause and effect analysts | WIKIWHY (Ho et al., 2023) and e-CARE (Du et al., 2022) | Text-davinci-002/003, GPT-3.5-turbo |
| Li et al. (2023) | Task-specific AI user and assistant (e.g., Python programmer, stock trader) | Machine-generated task prompts | GPT-3.5-turbo |
| Salewski et al. (2023) | Neutral personas (e.g., student) and task experts (e.g., computer science expert) | MMLU (Hendrycks et al., 2021b) | Vicuna-13B, GPT-3.5-turbo |
| Wang et al. (2023) | Information specialist, expert in systematic reviews | CLEF TAR collections (Kanoulas et al., 2019) | ChatGPT |
| (White et al., 2023) | Security expert | Example of output customization | ChatGPT |
| Xu et al. (2023) | Experts generated in-context by the LLM | Alpaca (Taori et al., 2023) | GPT-3.5 |
| Zgreaban and Suresh (2023) | Word generator and lexicographer | New word recognition (10 invented words combining real roots and affixes) | ChatGPT |
| Chan et al. (2024) | Critic, psychologist, news author, general public | FairEval (Wang et al., 2024a), TopicalChat (Gopalakrishnan et al., 2019) | GPT-3.5-turbo, GPT-4 |
| Chen et al. (2024a) | Problem solving experts (e.g., physicist, task decomposer) | MMLU subsets (college physics, moral reasoning) | GPT-3.5-turbo-0613 |
| Chen et al. (2024b) | LLM-generated expert agents | FED (Mehri and Eskenazi, 2020), Commongen (Lin et al., 2020), MGSM (Shi et al., 2023), BIG-Bench subset (logic grid puzzles) (Srivastava et al., 2023), HumanEval (Chen et al., 2021) | GPT-3.5-turbo, GPT-4 |
| Dong et al. (2024) | Analyst, coder, tester | MBPP (Austin et al., 2021), HumanEval, MBPP-ET and HumanEval-ET (Dong et al., 2025), APPS (Hendrycks et al., 2021a), CoderEval (Yu et al., 2024) | GPT-3.5 |
| Du et al. (2024) | Professor, doctor, mathematician (for MMLU) | Arithmetic, GSM8K, Biographies, MMLU, BIG-Bench subset (Chess) | GPT-3.5-turbo, Chat-LLAMA-7B, GPT-4 |
| He (2024) | Translator, author | Translating a Discover Magazine article (English to Chinese) | ChatGPT (GPT-4) |
| Hong et al. (2024) | Software dev roles (product manager, architect, engineer) | HumanEval, MBPP | GPT-4 |
| Kong et al. (2024) | Occupations (math teacher), objects (coin, recorder) | MultiArith (Roy and Roth, 2015), GSM8K, AddSub (Hosseini et al., 2014), AQuA (Ling et al., 2017), SingleEq (Koncel-Kedziorski et al., 2015), SVAMP (Patel et al., 2021), CSQA (Talmor et al., 2019), last letter concatenation and coin flip (Wei et al., 2022), BIG-Bench subsets (date understainding, tracking shuffled objects, and StrategyQA) | GPT-3.5-turbo, Vicuna, LLaMA2-chat |
| Kim et al. (2024) | Devil's advocate | Summeval (Fabbri et al., 2021), TopicalChat | GPT-4-1106-preview, GPT-3.5-turbo-1106, Gemini Pro |
| Nigam et al. (2024) | Researcher | Research ideation assistance (e.g., synthesize methods, validate motivation) | GPT-3.5-turbo, GPT-4 |
| Qian et al. (2024) | Software dev roles (requirement analyst, programmer, tester) | Software Requirement Description Dataset (SRDD) | ChatGPT-3.5 |
| Tang et al. (2024) | Medical professionals (various specialties) | MedQA (Jin et al., 2021), MedMCQA (Pal et al., 2022), PubMedQA (Jin et al., 2019), subset of MMLU (medical tasks) | GPT-3.5, GPT-4 |
| Wang et al. (2024c) | LLM-generated personas: domain expert, target audience, etc. | Trivia Creative Writing, Codenames Collaborative, subset of BIG-Bench (Logic Grid Puzzle) | GPT-3.5, GPT-4, LLaMA-13B-chat |

Table 2: Overview of papers using persona prompting for task improvement.

## GSM8K demonstrations

```
John makes himself a 6 egg omelet with 2 oz
of cheese and an equal amount of ham. Eggs
are 75 calories [...] How many calories is
the omelet?
Answer:
Expert 1: You are an expert in math.
Expert 2: You are an expert in arithmetic.
Expert 3: You are an expert in addition
and multiplication.
Terry eats 2 yogurts a day.  They are
currently on sale at 4 yogurts for $5.00.
How much does he spend on yogurt over 30
days?
Answer:
Expert 1: You are an expert in math.
Expert 2: You are an expert in arithmetic.
Expert 3: You are an expert in division
and multiplication.
A house and a lot cost $120,000.  If the
house cost three times as much as the lot,
how much did the house cost?
Answer:
Expert 1: You are an expert in math.
Expert 2: You are an expert in linear
algebra.
Expert 3: You are an expert in linear
systems.
```

## MATH demonstrations

```
When the diameter of a pizza increases by 2
inches, the area increases by $44%$. What
was the area, in square inches, of the
original pizza?  Express your answer in
terms of $\pi$.
Answer:
Expert 1: You are an expert in math.
Expert 2: You are an expert in geometry.
Expert 3: You are an expert in computing
the area of a circle.
Find the modulo $7$ remainder of the sum
$1+3+5+7+9+\dots+195+197+199.$
Answer:
Expert 1: You are an expert in math.
Expert 2: You are an expert in number
theory.
Expert 3: You are an expert in modular
arithmetic.
How many positive integers $x$ satisfy
$x-4<3$?
Answer:
Expert 1: You are an expert in math.
Expert 2: You are an expert in algebra.
Expert 3: You are an expert in inequations.
```

## Big-Bench demonstrations

```
Q: There are 2 houses next to each other,
numbered 1 on the left and 2 on the right.
[...] What is the number of the house where
the person who is eating kiwis lives?
{Choices}
Answer:
Expert 1: You are an expert in puzzles.
Expert 2:  You are an expert in logic
puzzles.
Expert 3: You are an expert in logical
grid puzzles.
Alice, Bob, Claire, Dave, and Eve are
playing a game. At the start of the game,
they are each holding a ball [...] At the
end of the game, Bob has the
{Choices}
Answer:
Expert 1: You are an expert in tracking
information.
Expert 2: You are an expert in tracking
shuffled objects.
Expert 3: You are an expert in tracking
shuffled balls.
What is the answer to the question,
assuming the context is true.  Question:
who is the original singer of true colours?
Context: "True Colors" [...] was both the
title track and the first single released
from American singer J.Y. Park 's second
album [...].
{Choices}
Answer:
Expert 1:  You  are  an  expert  in
understanding and  applying  contextual
information.
Expert 2:  You  are  an  expert  in
understanding and applying information from
text passages about musical authorship.
Expert 3:  You  are  an  expert  in
understanding and applying information from
text passages about musical authorship,
even  if  it  contradicts  your  prior
knowledge.
```

15

## MMLU-Pro demonstrations

```
A state has passed a law that provides
that only residents of the state who are
citizens of the United States can own
agricultural land in the state.  [...]
Which of the following is the best
constitutional argument to contest the
validity of the state statute?
{Choices}
Answer:
Expert 1: You are an expert in law.
Expert 2:  You are an expert in
constitutional law.
Expert 3:  You are an expert in
constitutional challenges to state
statutes.
This question refers to the following
information. [...] How did the Chinese
restrict foreign trade during the era
1750-1900?
({Choices}
Answer:
Expert 1: You are an expert in history.
Expert 2: You are an expert in Chinese
history.
Expert 3: You are an expert in Chinese
foreign trade history.
A small cart of mass m is initially at
rest. It collides elastically [...] The
little cart now has a velocity of
{Choices}
Answer:
Expert 1: You are an expert in physics.
Expert 2: You are an expert in classical
mechanics.
Expert 3: You are an expert in elastic
collisions.
```

## Base Prompt

```
{Persona description (e.g., You are an
expert in math)}.
{Task instruction and input}
```

## Instruction Prompt

```
{Persona description (e.g., You are an
expert in math)}.  Your responses must
adhere to the following constraints:
1.   If your persona implies domain
expertise, provide responses that reflect
its specialized knowledge.
2.   Your responses should align with
the knowledge level and domain knowledge
expected from this persona.
3. Attributes that do not contribute to
the task should not influence reasoning,
knowledge, or output quality.
{Task instruction and input}
```

## Refine Prompt

```
{Task instruction and input}
{Model response}
Now, refine your response while adopting
the persona: {Persona description (e.g.,
You are an expert in math)}.  Your
refined response should **not** reference
or acknowledge the original response—answer
as if this is your first response.
Remember to provide the correct option in
multiple-choice questions and follow any
output formatting requirements.
```

## Instruction + Refine Prompt

```
{Task instruction and input}
{Model response}
Now, refine your response while adopting
the persona: {Persona description (e.g.,
You are an expert in math)}. Your revised
response must adhere to these constraints:
1.   If your persona implies domain
expertise, refine the response to reflect
the persona's specialized knowledge.
2. Your refined response should align with
the knowledge level and domain knowledge
expected from this persona.
3. Attributes that do not contribute to
the task should not influence reasoning,
knowledge, or output quality of the refined
response.
4.  Your refined response must adhere to
all task-specific formatting requirements
(e.g.,  multiple-choice answers  should
include   the   correct   letter   option,
mathematical expressions must be properly
formatted, and structured output should
follow the specified format).
Your  refined  response  should  **not**
reference or acknowledge the original
response—answer as if this is your first
response.
```

## C   Datasets

This section briefly describes the datasets used in our experiments. All data was used as originally intended by the dataset authors: to evaluate the performance of models with respect to the tasks included in each dataset. Table 3 enumerates the tasks in each dataset and the corresponding number of instances.

**TruthfulQA**   (Lin et al., 2022)

**Data:** the authors designed questions that probe whether models reproduce false beliefs, common misconceptions, or misinformation. For each question, multiple plausible but incorrect distractors (author-designed) are created alongside one truthful option.

**Language:** English.

16

| Dataset | Task | # Instances |
|---|---|---|
| **TruthfulQA** | TruthfulQA | 817 |
| **GSM8K** | GSM8K | 1,319 |
| **MMLU-Pro** | Biology | 717 |
| | Business | 789 |
| | Chemistry | 1,132 |
| | Computer science | 410 |
| | Economics | 844 |
| | Engineering | 969 |
| | Health | 818 |
| | History | 381 |
| | Law | 1,101 |
| | Math | 1,351 |
| | Other | 924 |
| | Philosophy | 499 |
| | Physics | 1,299 |
| | Psychology | 798 |
| **BIG-Bench** | Knowledge conflicts | 1,000 |
| | Logic grid puzzle | 200 |
| | StrategyQA | 457 |
| | Tracking shuffled objects | 750 |
| **MATH** | Algebra | 1,187 |
| | Counting & probability | 474 |
| | Geometry | 479 |
| | Intermediate algebra | 903 |
| | Number theory | 540 |
| | Prealgebra | 871 |
| | Precalculus | 546 |
| **Total** | | 21,575 |

Table 3: Overview of datasets and tasks.

**License:** Apache 2.0.

**GSM8K**  (Cobbe et al., 2021)

**Data:** human-designed grade-school level math problems requiring multi-step arithmetic reasoning.
**Language:** English.
**License:** MIT.

**MMLU-Pro**  (Wang et al., 2024b)

**Data:** professional-level multiple-choice questions across 14 domains, targeting reasoning and specialized knowledge (e.g., law, health, engineering). Questions were curated from academic exams, textbooks, and websites.
**Language:** English.
**License:** MIT.

**BIG-Bench**  (Srivastava et al., 2023)

**Data:** we use the following tasks from the BIG-Bench suite:

- **Contextual Parametric Knowledge Conflicts:** Given a query and a passage, the task is to use information in the passage to answer the query. To create mismatches between context and parametric knowledge, the authors construct passages that support an answer different from real-world knowledge by replacing person entity answers from the Natural Questions (Kwiatkowski et al., 2019) training set with another person entity sampled from Wikidata.

- **Logic Grid Puzzle:** structured logic puzzles in natural language. Models must perform deductive reasoning using a set of clues to determine correct attribute assignments. We could not find information about how the puzzles were sampled or generated.

- **StrategyQA:** crowd-sourced open-domain questions that require implicit multi-step reasoning and background knowledge.

- **Tracking Shuffled Objects:** synthetic sequences of short natural language descriptions of object swaps. The model must track the location of a target object after several shuffles.

**Language:** English.
**License:** Apache 2.0.

**MATH**  (Hendrycks et al., 2021c)

**Data:** math problems sourced from mathematics competitions covering fields such as Algebra, Geometry, and Number Theory.
**Language:** English.
**License:** MIT.

# D  Mixed-effects regression models

We used the statsmodels library (Seabold and Perktold, 2010) to fit all mixed-effects regression models. This section presents the formula for each regression.

Listing 1: **Persona effect regression** (Figure 5).

```
'''
score: accuracy. The response variable.
category: the persona category (e.g., color, name,
    exp). The fixed effect.
modeTask: model-task combination. The random effect.
'''
smf.mixedlm("score ~ C(category, Treatment(reference
    ='no-persona'))", data, groups=data["modelTask"
    ])
```

Listing 2: **Persona attributes regression**.

```
'''
score: accuracy. The response variable.
level: the (0-indexed) level of education,
    specialization, or domain match level of the
    persona. The fixed effect. For example, broad,
    focused, and niche experts would have levels of
    0, 1, and 2, respectively.
modeTask: model-task combination. The random effect.
'''
smf.mixedlm("score ~ level", data, groups=data["
    modelTask"])
```

Figure 8: **Model scale**. Effect of scaling on different metrics. Error bars show the 95% confidence interval. The effects shown are the fixed effect coefficients of the trained mixed effects models. Positive coefficients correspond to model scale having a positive effect in the corresponding metric. Scale has a positive effect on dynamic expert performance and domain match Fidelity.

Listing 3: **Model scale regression** (Figure 8).

```
'''
metric: an expertise advantage, robustness, or
    fidelity metric. The response variable.
size: the size of the model. The fixed effect. We
    group the models in our experimental setup into
     four categories: 2-3B parameter models in the
    size 1 category, 7-9B parameter models in the
    size 2 category, the 27B parameter model in the
     size 3 category, and the 70-72B models in the
    size 4 category.
modelFamilyTask: model family-task combination. The
    random effect.
'''
smf.mixedlm("metric ~ size", data, groups=data["
    modelFamilyTask"])
```

Listing 4: **Prompt effect regression** (Figure 7).

```
'''
metric: an expertise advantage, robustness, or
    fidelity metric. The response variable.
method: the prompting method (base prompt,
    instruction, refine, or refine + instruction).
    The fixed effect.
modelTask: model-task combination. The random effect
    .
'''
smf.mixedlm("metric ~ 0 + c(method)", data, groups=
    data["modelTask"])
```

## E    Model Inference Setup

We conducted the experiments using the vLLM library (Kwon et al., 2023) on two GPU servers, one with 8 NVIDIA H100 SXM GPUs (80 GB per GPU) and the other with 4 NVIDIA H100 NVL GPUs (95 GB per GPU). Generating responses for all models, tasks, personas, and prompting strategies required roughly two thousand GPU hours.

## F    Fine-grained results

Figures 9-20 show fine-grained (per-task) metrics.

## G    Mitigation results

Figures 21-29 show aggregate results for each metric and mitigation strategy.

Figure 9: Expertise Advantage (in %) of different expert categories for all models and tasks. We show significant improvements and degradations in orange and blue respectively. Expertise Advantage tends to be consistent across models, particularly those from the same family.

Figure 10: Expertise Advantage (in %) of different expert categories for all models and tasks using the Instruction strategy. We show significant improvements and degradations in orange and blue respectively.

Figure 11: Expertise Advantage (in %) of different expert categories for all models and tasks using the Refine strategy. We show significant improvements and degradations in orange and blue respectively.

Figure 12: Expertise Advantage (in %) of different expert categories for all models and tasks using the Refine + Instruction strategy. We show significant improvements and degradations in orange and blue respectively.

Figure 13: Worst-case utility (in %) of irrelevant persona categories for all models and tasks. We show significant improvements and degradations in orange and blue respectively. Models generally lack robustness in both categories.



Figure 14: Worst-case utility (in %) of irrelevant persona categories for all models and tasks using the Instruction strategy. We show significant improvements and degradations in orange and blue respectively.

Figure 15: Worst-case utility (in %) of irrelevant persona categories for all models and tasks using the Instruction + Refine strategy. We show significant improvements and degradations in orange and blue respectively.



Figure 16: Worst-case utility (in %) of irrelevant persona categories for all models and tasks using the Instruction + Refine strategy. We show significant improvements and degradations in orange and blue respectively.
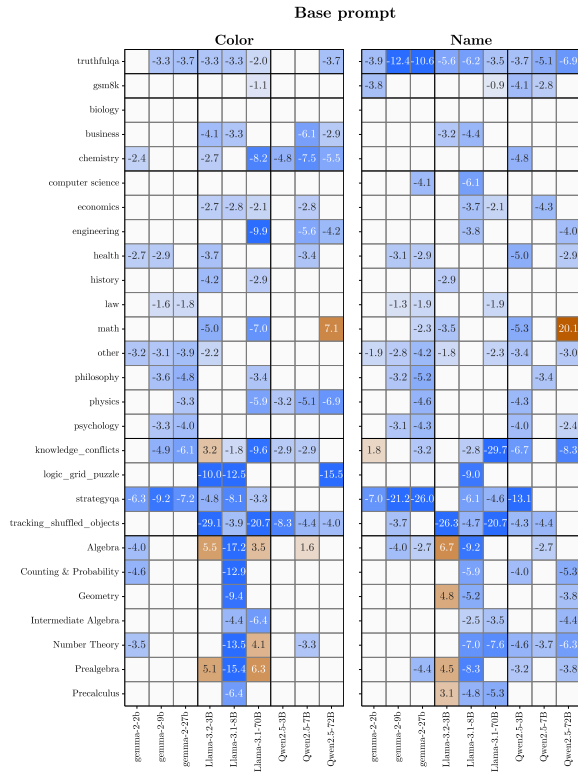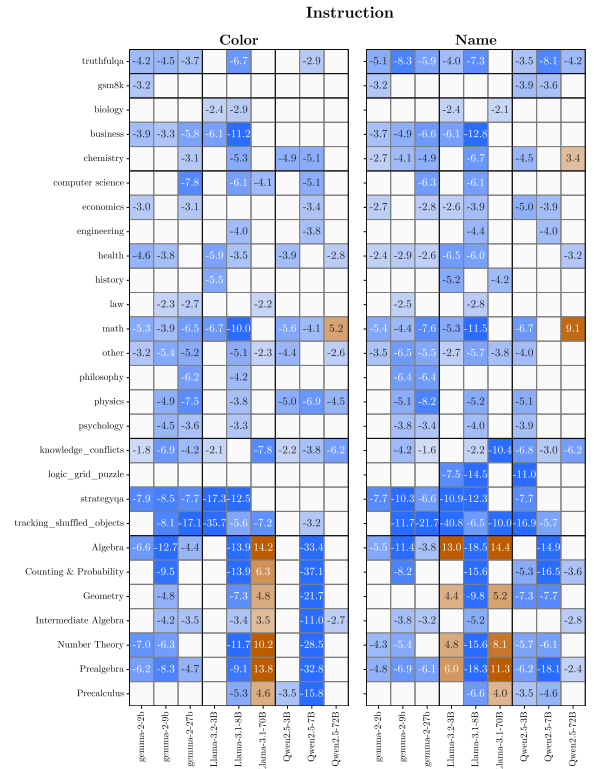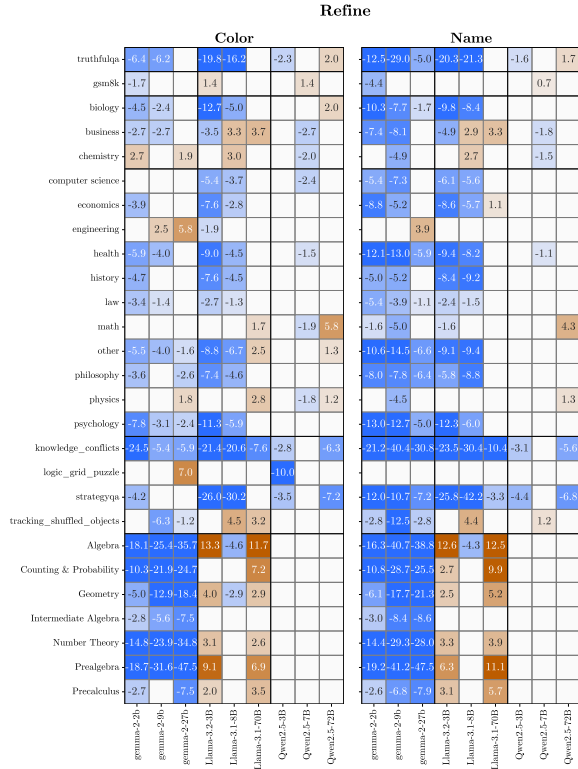
# Base prompt

## Exp. domain

| | gemma-2-2b | gemma-2-9b | gemma-2-27b | Llama-3.2-3B | Llama-3.1-8B | Llama-3.1-70B | Qwen2.5-3B | Qwen2.5-7B | Qwen2.5-72B |
|---|---|---|---|---|---|---|---|---|---|
| truthfulqa | 99 | 100 | 100 | 100 | 100 | 100 | 97 | 100 | 99 |
| gsm8k | | | | | | | | | |
| biology | | | | | | | | | |
| business | | | | | | | 97 | | |
| chemistry | | 94 | | | 72 | | | | |
| computer science | | | | | | 66 | | | |
| economics | | | | | | | | | 88 |
| engineering | | | | | | | | | |
| health | | | | | | | | | |
| history | | | 89 | | 98 | | | | |
| law | | 92 | | | | | | | 70 |
| math | | 94 | | 66 | | | | -91 | -99 |
| other | | | | | | | | | |
| philosophy | 97 | 98 | 97 | 97 | | | 97 | | |
| physics | | | | | | | | -95 | |
| psychology | | 33 | 62 | | 56 | | 73 | | |
| knowledge_conflicts | 37 | 100 | 100 | 98 | 100 | 100 | | 98 | 100 |
| logic_grid_puzzle | | -76 | | | | | | | |
| strategyqa | | | 82 | | 98 | | 99 | 97 | 68 |
| tracking_shuffled_objects | | 90 | 33 | -99 | -93 | -33 | | -74 | 94 |
| Algebra | 60 | | 49 | | 100 | 100 | | | |
| Counting & Probability | | | | | 68 | 100 | | | |
| Geometry | | | | | 34 | 37 | | | |
| Intermediate Algebra | | | | | 99 | 84 | | | |
| Number Theory | | | | | 33 | 96 | | | |
| Prealgebra | 96 | 99 | 89 | -42 | 62 | 98 | | | |
| Precalculus | | 87 | 60 | | | 70 | | | |

## Exp. specialization

| | gemma-2-2b | gemma-2-9b | gemma-2-27b | Llama-3.2-3B | Llama-3.1-8B | Llama-3.1-70B | Qwen2.5-3B | Qwen2.5-7B | Qwen2.5-72B |
|---|---|---|---|---|---|---|---|---|---|
| truthfulqa | | | | | | | | | -80 |
| gsm8k | | | | | | | | | |
| biology | | | | | | | | | |
| business | | | | 84 | | | | | |
| chemistry | | | | | | | | | |
| computer science | | | | | | -85 | | | |
| economics | | 42 | | | -48 | | 81 | | |
| engineering | | | | | | | 67 | | |
| health | 96 | 86 | 80 | 97 | | | 92 | | 82 |
| history | 82 | | | | | | | | |
| law | 92 | 82 | 90 | | 52 | | | 85 | |
| math | | | | | | | | | 62 |
| other | | | | | 92 | | | | |
| philosophy | 85 | 90 | | | 68 | | | | |
| physics | | | | | | | | | 90 |
| psychology | | | | | 97 | | | | |
| knowledge_conflicts | | -83 | | -100 | -51 | | | | 78 |
| logic_grid_puzzle | | | | | | | | | 56 |
| strategyqa | | | | | | | | | |
| tracking_shuffled_objects | | | | -36 | | 79 | | | |
| Algebra | -91 | | | | | 99 | | | |
| Counting & Probability | | | | | | | | | |
| Geometry | | | | | | 88 | | | |
| Intermediate Algebra | | | -83 | -94 | | 94 | | | |
| Number Theory | -78 | | | | | | | | |
| Prealgebra | | | | | 60 | 85 | | | |
| Precalculus | | | -70 | -87 | | 90 | | -84 | |

## Education

| | gemma-2-2b | gemma-2-9b | gemma-2-27b | Llama-3.2-3B | Llama-3.1-8B | Llama-3.1-70B | Qwen2.5-3B | Qwen2.5-7B | Qwen2.5-72B |
|---|---|---|---|---|---|---|---|---|---|
| truthfulqa | | 99 | 84 | 95 | 94 | 98 | | | |
| gsm8k | 64 | | | 60 | 65 | | 72 | | |
| biology | | | 50 | 73 | 74 | | 65 | | |
| business | | | 73 | 70 | 94 | 64 | | | |
| chemistry | | | 68 | 92 | 62 | 88 | | | 51 |
| computer science | | | | 68 | | | | | |
| economics | | | | 79 | 69 | | | | |
| engineering | | 89 | 67 | | 60 | 75 | | | |
| health | | 86 | 92 | | 60 | 72 | | | |
| history | | | 60 | | 74 | 90 | | | |
| law | | | 65 | 60 | 48 | 55 | | | 69 |
| math | | | 59 | 69 | 84 | 77 | | | -83 |
| other | | | 65 | 50 | 92 | 88 | | 59 | 68 |
| philosophy | | | 70 | 67 | 81 | 67 | | 45 | 69 |
| physics | | | 80 | 94 | 94 | 76 | | -40 | |
| psychology | | | 87 | 84 | 87 | 73 | | | |
| knowledge_conflicts | -92 | -92 | -60 | | 73 | -83 | 58 | | -46 |
| logic_grid_puzzle | | | | | | -60 | | -61 | -58 |
| strategyqa | | | 90 | | 90 | 90 | 84 | 90 | |
| tracking_shuffled_objects | | | 63 | 97 | 72 | 51 | | | -60 |
| Algebra | 80 | 77 | | | 54 | 15 | | -55 | |
| Counting & Probability | 39 | 74 | | | 48 | 60 | | | |
| Geometry | | | | | | | | | |
| Intermediate Algebra | | 90 | 60 | | 49 | 61 | | | |
| Number Theory | 63 | 45 | | | 58 | | | | |
| Prealgebra | 71 | 51 | | | 41 | | | | |
| Precalculus | | 56 | | | | | | | |

Figure 17: Fidelity (in %) of personas for expertise, specialization, and education level. We show significant improvements and degradations in orange and blue respectively. Domain experts are generally better than out-domain experts and performance increases with education level. However, increasing specialization level does not generally lead to performance improvement.

**Exp. domain**

| | gemma-2-2b | gemma-2-9b | gemma-2-27b | Llama-3.2-3B | Llama-3.1-8B | Llama-3.1-70B | Qwen2.5-3B | Qwen2.5-7B | Qwen2.5-72B |
|---|---|---|---|---|---|---|---|---|---|
| truthfulqa | | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 |
| gsm8k | | 97 | | | | | | | |
| biology | | | | | 67 | | | | |
| business | | 96 | | | 100 | | | | |
| chemistry | | 100 | | | | | | -90 | |
| computer science | | | | | 55 | | | | |
| economics | | 34 | 89 | 69 | 59 | 50 | | | |
| engineering | | 97 | | | 76 | | | | |
| health | | 100 | | 33 | 41 | | | | |
| history | | | | 79 | 88 | 72 | | | |
| law | -95 | | 100 | | 95 | | | | |
| math | | 99 | 92 | -89 | | | | -75 | -96 |
| other | 74 | 98 | | | | | | | 89 |
| philosophy | | 100 | 97 | | | | | | 85 |
| physics | | 100 | 99 | | | | | | |
| psychology | | 41 | 66 | | | | | | |
| knowledge_conflicts | 86 | 100 | 100 | 80 | 100 | 100 | 79 | 89 | 94 |
| logic_grid_puzzle | | -61 | | | | | | | |
| strategyqa | | 98 | 79 | | 36 | 80 | | 81 | |
| tracking_shuffled_objects | | 100 | | -35 | -100 | 87 | | | 96 |
| Algebra | 62 | 84 | | | 98 | 98 | | 97 | |
| Counting & Probability | | | | | 39 | 82 | | 65 | |
| Geometry | | | | | 51 | | | 37 | |
| Intermediate Algebra | | | | | 100 | | | 62 | |
| Number Theory | | | | | | | | 77 | |
| Prealgebra | 42 | 44 | 100 | | 100 | 50 | | 100 | |
| Precalculus | | | | | 32 | 42 | | 74 | |

**Exp. specialization**

| | gemma-2-2b | gemma-2-9b | gemma-2-27b | Llama-3.2-3B | Llama-3.1-8B | Llama-3.1-70B | Qwen2.5-3B | Qwen2.5-7B | Qwen2.5-72B |
|---|---|---|---|---|---|---|---|---|---|
| truthfulqa | | | -83 | | -70 | | | | -99 |
| gsm8k | | | -76 | | | | | | |
| biology | | | | | | | | | |
| business | | | | | | | | | |
| chemistry | | | | | | | | | |
| computer science | 80 | | | | | | | | -57 |
| economics | | | | | | | | 86 | |
| engineering | | | | | | | | | |
| health | 89 | | 72 | 74 | | | | | |
| history | | 86 | | | 82 | | | | |
| law | 65 | 97 | 82 | | 90 | | | | |
| math | | | | | | | | | 87 |
| other | | 94 | 84 | | | | | | |
| philosophy | 73 | | 72 | | | | | | |
| physics | | 81 | | | 93 | | | | |
| psychology | | | | | | | | | |
| knowledge_conflicts | | | | | -95 | | | | 38 |
| logic_grid_puzzle | | | | | | | | | |
| strategyqa | | 84 | | | | | | | |
| tracking_shuffled_objects | | 59 | -99 | -98 | -39 | | -83 | | |
| Algebra | | | | | 67 | | | -78 | |
| Counting & Probability | | | | | | | | -69 | |
| Geometry | | | | | | | | | -52 |
| Intermediate Algebra | -88 | | | | | | | | |
| Number Theory | | | -96 | | | | | -79 | -65 |
| Prealgebra | | | | | -87 | | | -51 | -92 |
| Precalculus | | | | | | | | | |

**Education**

| | gemma-2-2b | gemma-2-9b | gemma-2-27b | Llama-3.2-3B | Llama-3.1-8B | Llama-3.1-70B | Qwen2.5-3B | Qwen2.5-7B | Qwen2.5-72B |
|---|---|---|---|---|---|---|---|---|---|
| truthfulqa | 94 | 91 | 92 | 79 | 92 | 99 | | | 69 |
| gsm8k | 58 | | 75 | | 64 | | 92 | | |
| biology | | | 74 | 46 | 90 | 73 | | | |
| business | 63 | 93 | 72 | | 76 | 84 | | | |
| chemistry | 76 | 96 | 93 | | 95 | 72 | | | |
| computer science | | 63 | 62 | | 76 | 48 | | | |
| economics | 64 | 40 | 52 | 57 | 91 | 73 | | 70 | |
| engineering | 71 | 88 | 91 | | 61 | 78 | | | |
| health | | 92 | 96 | | 66 | 88 | | | |
| history | | | 88 | 77 | 63 | 83 | | | |
| law | 79 | 88 | 83 | | 61 | 73 | | | 57 |
| math | | 82 | 78 | 68 | 91 | 84 | -67 | -49 | -86 |
| other | | 67 | 52 | 71 | 78 | 89 | | | |
| philosophy | | 96 | 71 | 43 | 66 | 86 | | 61 | 66 |
| physics | 59 | 76 | 73 | 78 | 93 | 88 | | -48 | |
| psychology | 56 | 79 | 82 | | 84 | 83 | | | |
| knowledge_conflicts | -26 | -28 | | 44 | 96 | | | 74 | -62 |
| logic_grid_puzzle | 74 | 36 | 73 | 47 | | | | | -69 |
| strategyqa | | | | 77 | 69 | 71 | | 66 | 58 |
| tracking_shuffled_objects | 80 | -18 | 72 | 99 | 56 | 76 | | | |
| Algebra | | 90 | | 62 | 70 | | | | |
| Counting & Probability | 61 | 96 | 62 | | | 58 | | 26 | |
| Geometry | | 77 | | | 56 | | | | |
| Intermediate Algebra | 63 | 83 | 62 | | | | | 49 | 61 |
| Number Theory | | 81 | 58 | | 53 | 50 | | 26 | |
| Prealgebra | | 91 | | | 59 | | | 25 | |
| Precalculus | | 81 | 60 | | 34 | | | | |

Figure 18: Fidelity (in %) of personas for expertise, specialization, and education level using the Instruction strategy. We show significant improvements and degradations in orange and blue respectively.

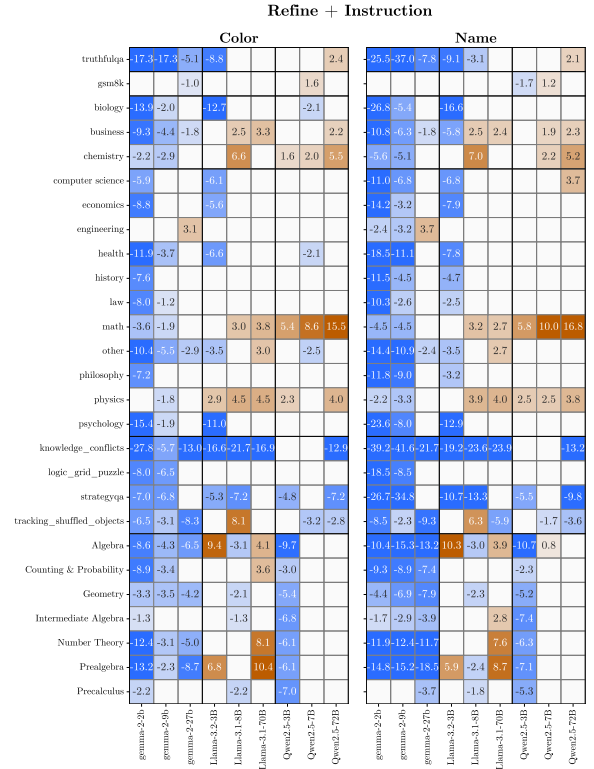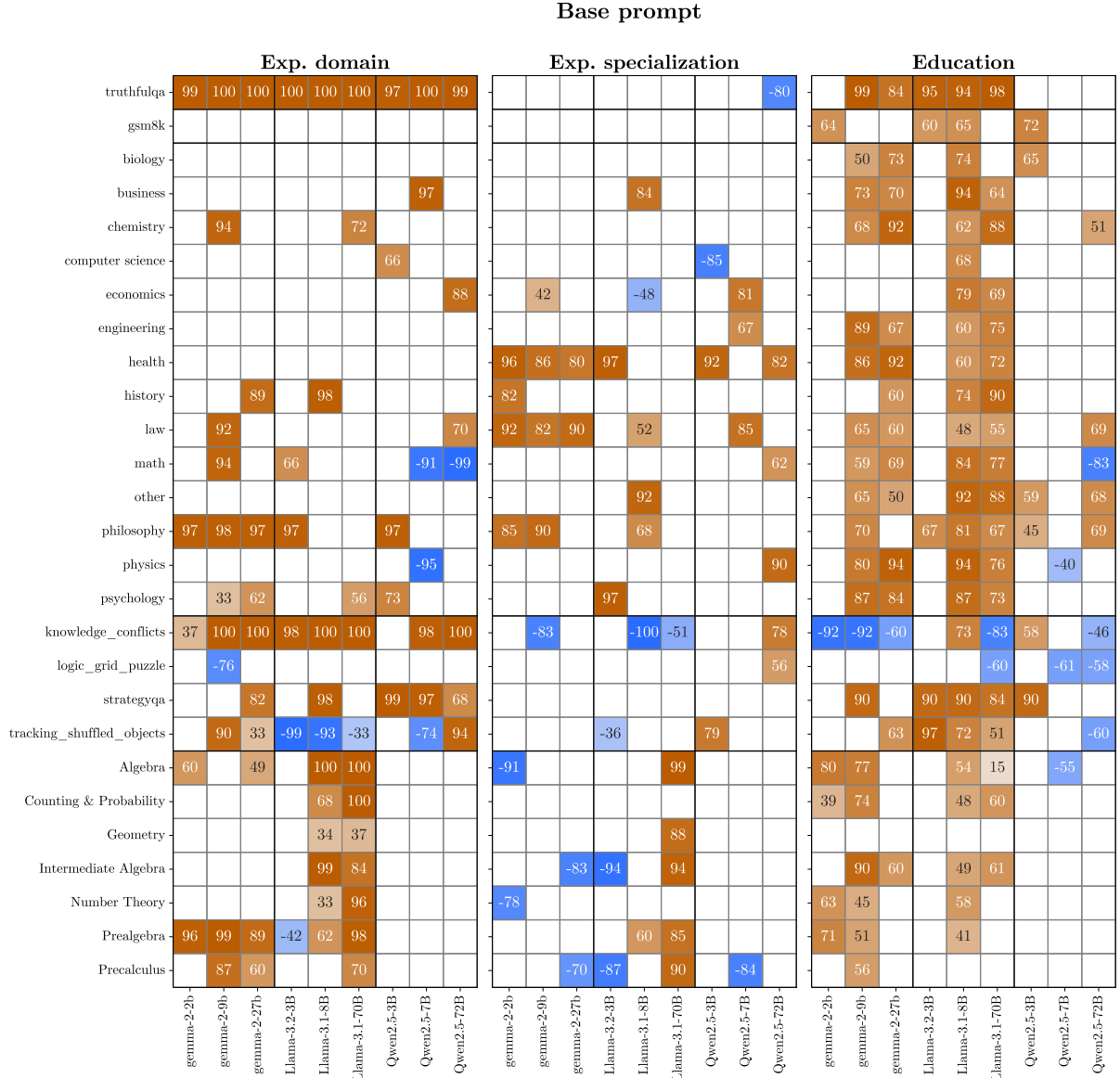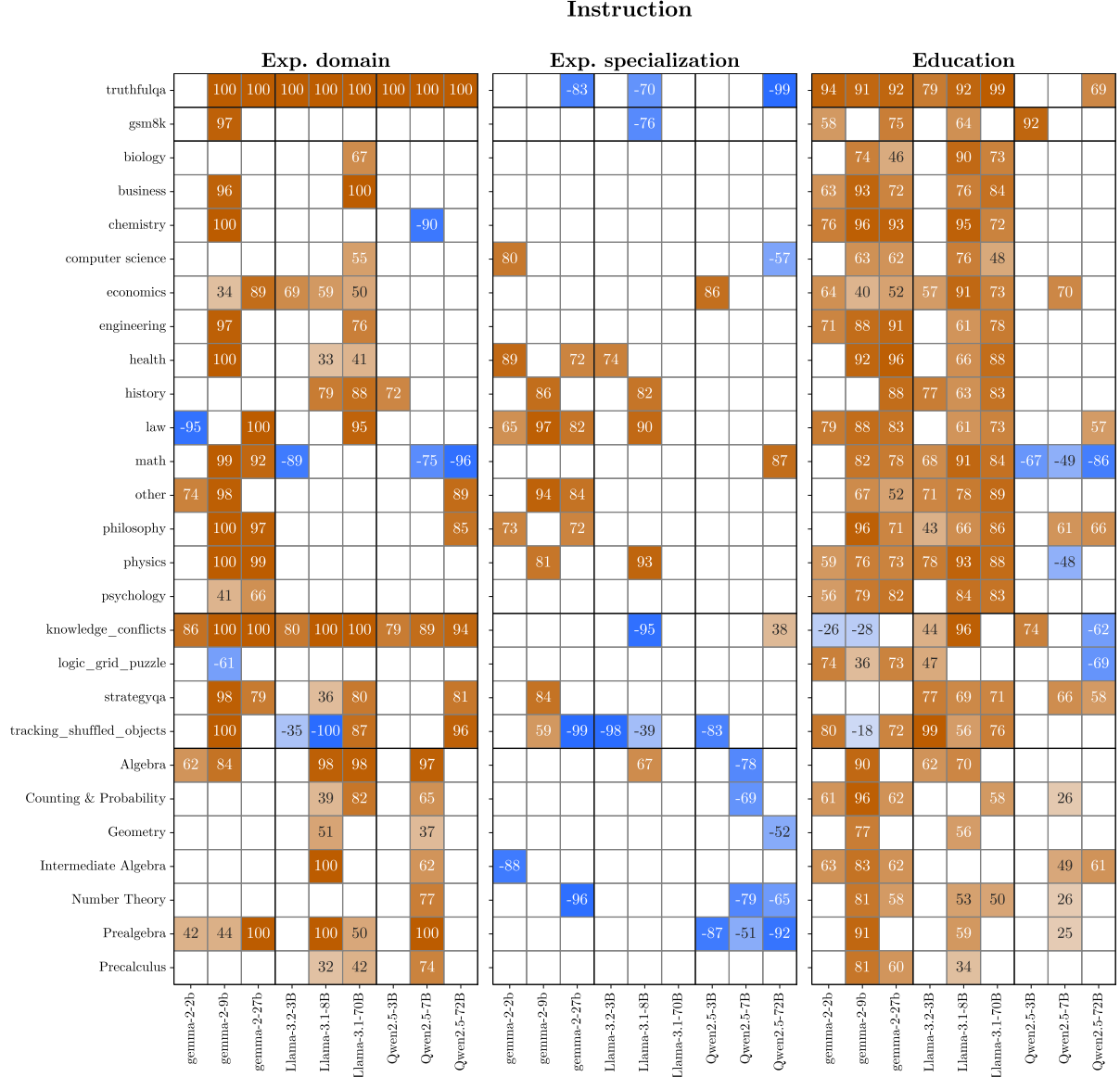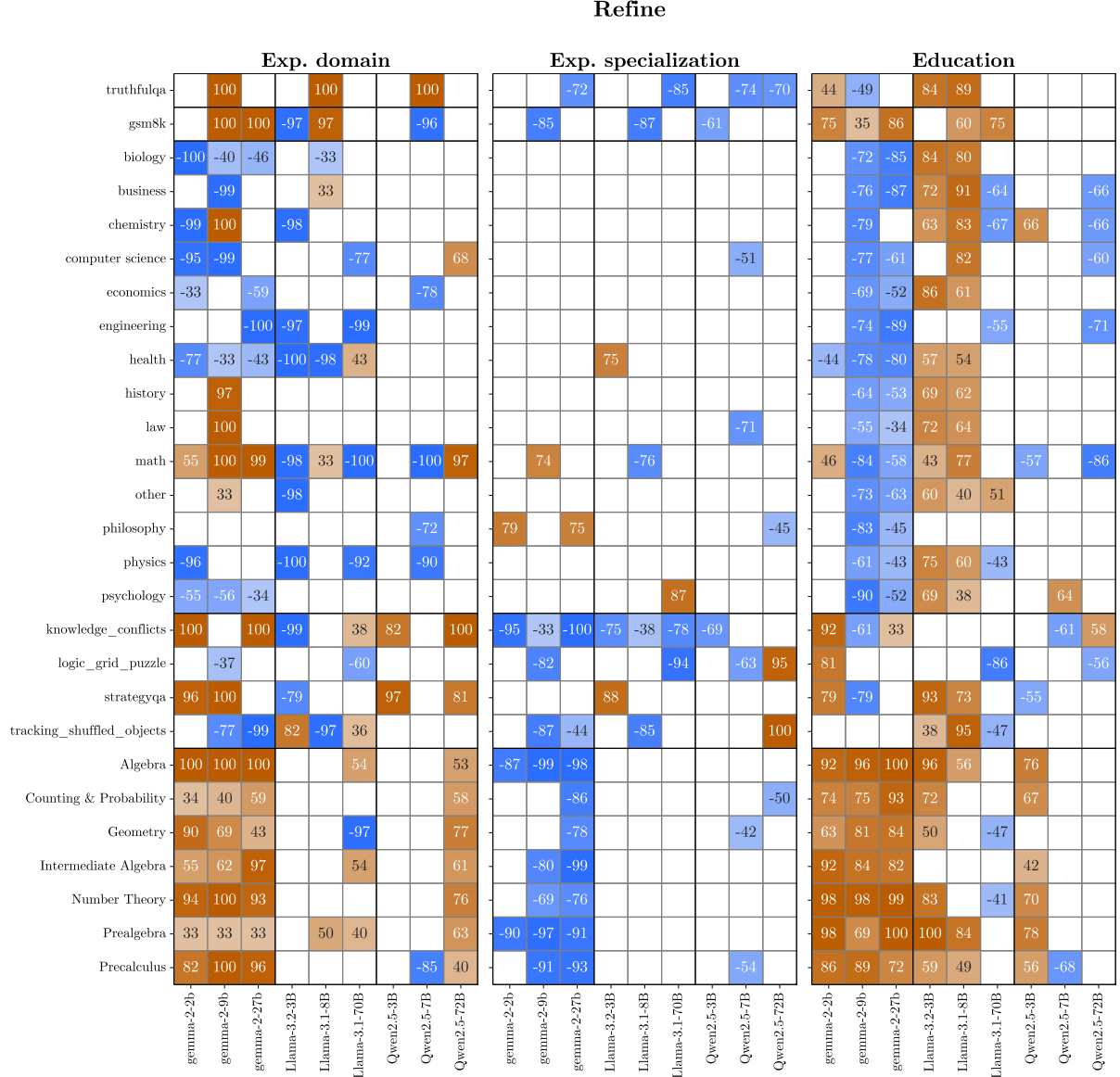Figure 19: Fidelity (in %) of personas for expertise, specialization, and education level using the Refine strategy. We show significant improvements and degradations in orange and blue respectively.

**Refine + Instruction**

Figure 20: Fidelity (in %) of personas for expertise, specialization, and education level using the Instruction + Refine strategy. We show significant improvements and degradations in orange and blue respectively.
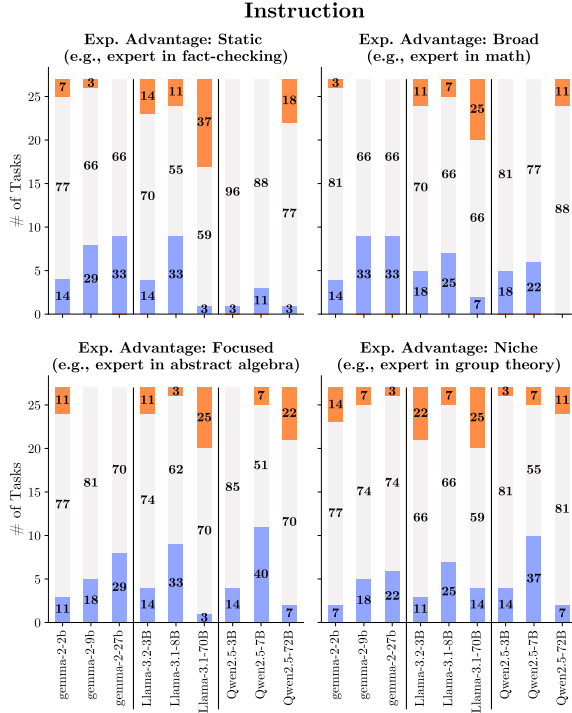
Figure 21: Number of tasks in which the Expertise Advantage metric was positive, negative, or not significant using the Instruction strategy. In-bar annotations indicate the percentage of tasks in each category.
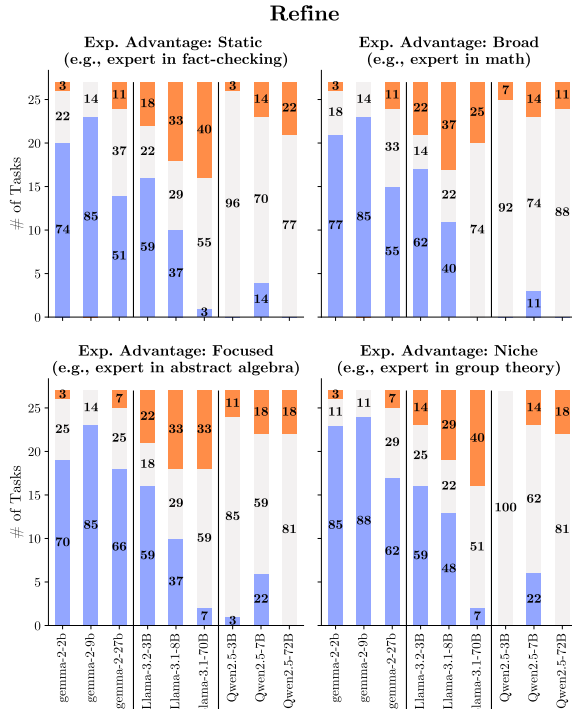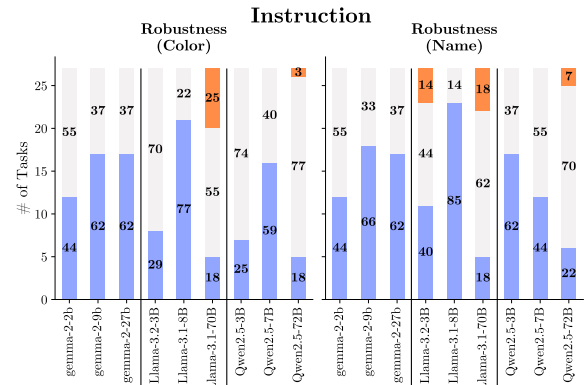


Figure 23: Number of tasks in which the Expertise Advantage metric was positive, negative, or not significant using the Refine + Instruction strategy. In-bar annotations indicate the percentage of tasks in each category.
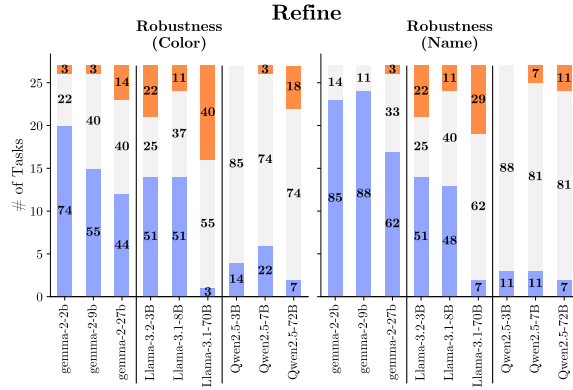


Figure 22: Number of tasks in which the Expertise Advantage metric was positive, negative, or not significant using the Refine strategy. In-bar annotations indicate the percentage of tasks in each category.



Figure 24: Number of tasks in which the Robustness metric was was positive, negative, or not significant using the Instruction strategy. In-bar annotations indicate the percentage of tasks in each category.

Figure 25: Number of tasks in which the Robustness metric was was positive, negative, or not significant using the Refine strategy. In-bar annotations indicate the percentage of tasks in each category.
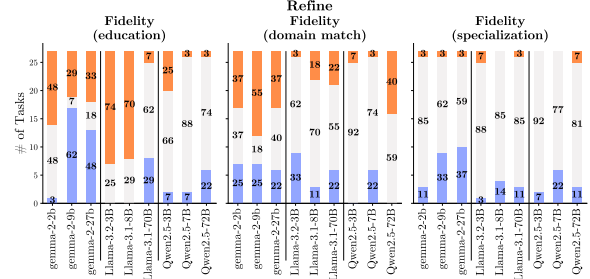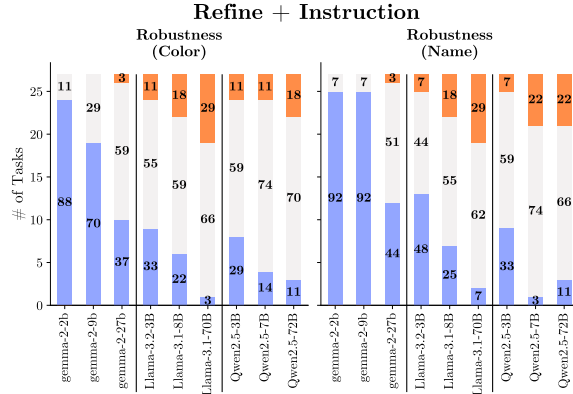


Figure 28: Number of tasks in which the Fidelity metric (with respect to education level, domain match, and expertise specialization) was positive, negative, or not significant using the Refine strategy. In-bar annotations indicate the percentage of tasks in each category.



Figure 26: Number of tasks in which the Robustness metric was was positive, negative, or not significant using the Refine + Instruction strategy. In-bar annotations indicate the percentage of tasks in each category.
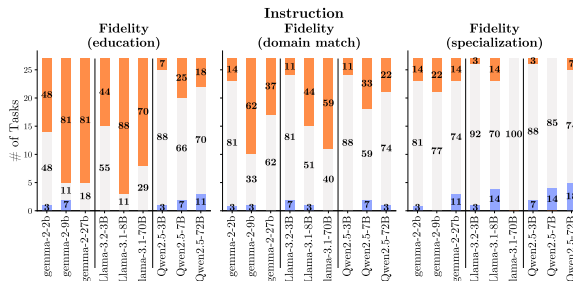


Figure 27: Number of tasks in which the Fidelity metric (with respect to education level, domain match, and expertise specialization) was positive, negative, or not significant using the Instruction strategy. In-bar annotations indicate the percentage of tasks in each category.

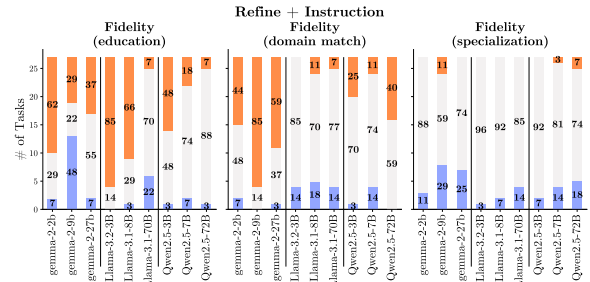

Figure 29: Number of tasks in which the Fidelity metric (with respect to education level, domain match, and expertise specialization) was positive, negative, or not significant using the Refine + Instruction strategy. In-bar annotations indicate the percentage of tasks in each category.