

BENCHMARK FOR ASSESSING OLFACTORY PERCEPTION OF LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Here we introduce the Olfactory Perception (OP) benchmark, designed to assess the capability of large language models (LLMs) to reason about smell. The benchmark contains 1,010 questions across eight task categories spanning odor classification, odor primary descriptor identification, intensity and pleasantness judgments, multi-descriptor prediction, mixture similarity, olfactory receptor activation, and smell identification from real-world odor sources. Each question is presented in two prompt formats, compound names and isomeric SMILES, to evaluate the effect of molecular representations. Evaluating 21 model configurations across major model families, we find that compound-name prompts consistently outperform isomeric SMILES, with gains ranging from +2.4 to +18.9 percentage points (mean \approx +7 points), suggesting current LLMs access olfactory knowledge primarily through lexical associations rather than structural molecular reasoning. The best-performing model reaches 64.4% overall accuracy, which highlights both emerging capabilities and substantial remaining gaps in olfactory reasoning. We further evaluate a subset of the OP across 21 languages and find that aggregating predictions across languages improves olfactory prediction, with AUROC = 0.86 for the best performing language ensemble model. LLMs should be able to handle olfactory and not just visual or aural information.

1 INTRODUCTION

The sense of smell occupies a unique position among human perceptual modalities (Mainland et al., 2015; Secundo et al., 2014). Unlike vision, where wavelengths map predictably to colors (Stockman & Sharpe, 2000), or audition, where frequencies correspond to pitch (Oxenham et al., 2011), olfaction operates through a complex interplay between molecular structure and receptor biology that remains incompletely understood (Mainland et al., 2015; Buck & Axel, 1991; Bushdid et al., 2014). Humans can distinguish over a trillion distinct odors Bushdid et al. (2014); err (2024), though estimates are still debated (Gerkin & Castro, 2015), yet predicting how a molecule will smell from its chemical structure alone has long eluded computational approaches (Keller et al., 2017). This structure-odor relationship represents one of the most challenging frontiers in sensory science Keller et al. (2017); Lee et al. (2023); Secundo et al. (2014).

Large language models have demonstrated remarkable capabilities across diverse domains, from mathematical reasoning to code generation Brown et al. (2020); Chen et al. (2021). Recent work has shown that these models, particularly those trained with explicit reasoning objectives (Han et al., 2024), can perform sophisticated tasks in chemistry, interpreting SMILES (Simplified Molecular Input Line Entry System) (Weininger, 1988) strings, predicting reaction products, and even elucidating molecular structures from NMR spectra Runcie et al. (2025); Mirza et al. (2025). These advances raise a natural question: can LLMs bridge the gap between molecular structure and sensory perception? Specifically, do these models possess knowledge about how molecules smell?

This question carries both scientific and practical significance. LLMs capable of olfactory reasoning could accelerate fragrance design Mao et al. (2018), flavor development Ge et al. (2025), and the identification of malodorous contaminants Bartsch et al. (2016). Yet despite growing interest in evaluating LLM sensory alignment across modalities including color, auditory perception, and taste Marjieh et al. (2024), olfaction has remained notably absent from systematic evaluation. Prior work has found limited human-AI alignment for smell Zhong et al. (2024), shown that LLMs capture

054 olfactory-semantic similarity but not factual accuracy Kurfalı et al. (2025), and demonstrated that
055 specialized graph neural networks can achieve human-level odor prediction Keller et al. (2017); Lee
056 et al. (2023). However, no existing work systematically evaluates whether general-purpose LLMs
057 possess factual knowledge about olfactory properties through structured question-answering.

058 We address this gap by introducing the Olfactory Perception (OP) benchmark, comprising 1,010
059 questions across eight task categories that span the full complexity of olfactory perception: from basic
060 odor detection Mayhew et al. (2022) and classification to intensity and pleasantness judgments Keller
061 et al. (2017), from single-molecule descriptor identification International Fragrance Association
062 (2020) to mixture similarity Bushdid et al. (2014); Snitz et al. (2013); Ravia et al. (2020); Satarifard
063 et al. (2025), and from semantic odor labeling Lee et al. (2023) to receptor activation Lalis et al.
064 (2024). Ground-truth answers are derived from established datasets in olfactory science. Motivated
065 by cross-linguistic diversity in odor language Majid & Burenhult (2014); Majid et al. (2018); Majid
066 (2021), we additionally evaluate a subset across 21 languages and report a cross-lingual majority-vote
067 setting.

068 A distinctive feature of our approach is the dual-prompting strategy: each question is presented
069 using both isomeric SMILES notation and common compound names. Because odor perception is
070 stereospecific (e.g., (R)-carvone smells of spearmint while (S)-carvone smells of caraway), this design
071 enables direct comparison of how molecular representation format affects olfactory reasoning Runcie
072 et al. (2025), assessing whether models reason about molecular properties or merely retrieve lexical
073 associations. Our evaluation encompasses 21 model configurations across six providers, including
074 OpenAI’s GPT and o3 series, Google’s Gemini, Anthropic’s Claude, Meta’s Llama, xAI’s Grok, and
075 DeepSeek. We systematically vary reasoning budgets where applicable, enabling analysis of how
076 extended deliberation affects olfactory task performance, an approach that has revealed substantial
077 performance differences in chemistry tasks Runcie et al. (2025); Mirza et al. (2025). Our findings
078 reveal that while current models achieve moderate success on certain tasks, substantial gaps remain,
079 particularly for questions requiring genuine molecular reasoning rather than factual recall about
080 well-known compounds.

081 The contributions of this work are threefold. First, we introduce a comprehensive benchmark for
082 evaluating LLMs on olfactory reasoning, comprising 1,010 questions across eight task categories
083 grounded in peer-reviewed olfactory science. Second, we establish baseline performance across
084 state-of-the-art models, identifying both emerging capabilities and systematic failures. Third, via
085 dual-prompting (compound names vs. isomeric SMILES), we provide new insight into how molecular
086 representation shapes model performance, helping to distinguish structural reasoning from lexical
087 association.

089 2 RELATED WORK

091 Recent benchmarks have evaluated LLM chemical reasoning, including ChemBench Mirza et al.
092 (2025), ChemIQ Runcie et al. (2025), LlaSMol Yu et al. (2024), GPQA Rein et al. (2024), MMLU
093 Hendrycks et al. (2021), and a comprehensive eight-task evaluation Guo et al. (2023), but none test
094 whether LLMs can reason about how molecules smell. More broadly, LLMs show substantial but
095 incomplete alignment with human sensory judgments across modalities such as pitch, color, and taste
096 Marjeh et al. (2024); Xu et al. (2025); Fukushima et al. (2025), with performance degrading for
097 modalities less frequently described in language. Olfaction, notably absent from these evaluations,
098 represents a natural test case.

099 Initial work on LLMs and olfaction has found limited human-AI alignment: SNIFF AI Zhong et al.
100 (2024) achieved only 27.5% success in scent identification from descriptions, and Kurfalı et al. (2025)
101 showed that GPT-4o captures olfactory-semantic similarity but evaluated similarity judgments rather
102 than factual accuracy. Related multimodal work has explored image-text olfactory matching Kurfalı
103 et al. (2023); Esteban-Romero et al. (2025). Meanwhile, specialized ML systems purpose-built for
104 olfaction have achieved strong performance: the Principal Odor Map (POM) Lee et al. (2023) reaches
105 AUROC = 0.89 for multi-label descriptor prediction on the GS-LF dataset (the same 138-descriptor
106 vocabulary used in our RATA task), Dense Sense Saha et al. (2025) reaches AUROC = 0.875, Mol-
107 PECO Zhang et al. (2023) achieves AUROC of 0.813, and POMMIX Tom et al. (2025) extends
the POM to mixture similarity (Pearson $\rho = 0.779$). Two DREAM Challenges Keller et al. (2017);

Satarifard et al. (2025) provided foundational benchmarks for these systems. However, none of this work addresses whether general-purpose LLMs possess latent olfactory knowledge.

Our OP benchmark fills this gap as, to our knowledge, the first structured factual evaluation of LLM olfactory reasoning, comprising 1,010 questions across eight task categories with ground-truth answers derived from established olfactory science datasets (Mayhew et al., 2022; Keller et al., 2017; Lee et al., 2023; Lalis et al., 2024; Snitz et al., 2013; Bushdid et al., 2014; Ravia et al., 2020; International Fragrance Association, 2020; Kreissl et al., 2022). A dual-prompting strategy (isomeric SMILES vs. compound names) enables direct comparison of how molecular representation affects olfactory reasoning. A detailed discussion and side-by-side comparison with prior work is provided in Appendix B.

3 OLFACTORY PERCEPTION BENCHMARK

We introduce a unified benchmark of 1,010 olfaction questions spanning odor detectability, semantic odor description, perceptual judgments, mixture similarity, receptor activation, and smell identification test from mixtures. Figure 1 provides an overview of the benchmark; a detailed summary of question categories, data sources, and response formats is provided in Appendix A. Each item is presented in two equivalent formats: (i) isomeric SMILES (prompt 1), and (ii) compound names (prompt 2), to separate structure-based reasoning from name priors. Tasks are multiple choice from constrained lists of options to enable consistent automatic scoring across models.

Odor Classification (OC). Odor classification is a binary task that asks whether a given molecule is *Odorous* or *Odorless* (175 questions, 50% Odorous). This task assesses basic odor detectability prediction from chemical identity of a molecule. Molecules were obtained from a previously curated dataset Mayhew et al. (2022) and molecular weight was constrained to be ≤ 350.0 (so as to be in the smellable range); items were presented to LLMs in random order. We consider odor classification to be a simple task.

Odor Primary Descriptor (OPD). Odor primary descriptor is a multi-class task where the model selects the single descriptor for a molecule from a provided list of four options (175 questions). This task evaluates whether models can map chemicals to main semantic odor categories. The molecules were obtained from the 2020 version of the IFRA fragrance ingredient glossary (FIG) International Fragrance Association (2020). Besides the primary descriptor, FIG provides secondary and tertiary descriptors; we excluded these descriptors from the list of distractors to obtain a more robust evaluation. Selection of molecules was done on a stratified set of descriptors to represent less dominant descriptors, and a total of 29 descriptors are present among questions. We consider odor primary descriptor a simple task.

Odor Intensity and Pleasantness (OIn and OPI). To evaluate model assessment of two olfactory dimensions (*i.e.*, intensity and pleasantness), we use two paired-comparison tasks, where the model chooses which of two molecules is more intense or more pleasant (175 molecular pairs for intensity, and 175 molecular pairs for pleasantness comparisons). Ground truth data are obtained from prior work Keller et al. (2017) which provides mean human subjective ratings on a 0–100 scale. Molecular pairs were selected to be above and below median with a minimum score difference of 10. Furthermore, prompts ask the model to rate the intensity and pleasantness on a 0–100 scale. We consider odor intensity and pleasantness categories to be simple tasks.

Rate-All-That-Apply (RATA). To assess more complex olfactory perceptual capability, we employ a multilabel semantic profiling task (100 questions), where, given a molecule and a descriptor lexicon (138 odor descriptors), the model selects all descriptors that apply in describing the odor of the molecule. We selected 100 molecules from an integrated dataset of the Good Scents and Leffingwell Associates (GS-LF) Lee et al. (2023). We constrained our selection to molecules with 2 to 5 answers from 138 descriptors, with 25 questions for 2, 3, 4, and 5 descriptors, respectively. We consider RATA categorization as an intermediate-difficulty task.

Odor Similarity of Mixtures (OS). In odor similarity of mixtures, we evaluate a categorical mixture-perception task (100 questions) where the model compares two mixtures (each a set of 2–10 molecules) and predicts an ordinal similarity label (*e.g.*, from strongly similar to strongly dissimilar). We selected 100 mixture pairs from various datasets Snitz et al. (2013); Bushdid et al. (2014); Ravia et al. (2020) and after standardization of perceptual distance, similarity values were categorized into

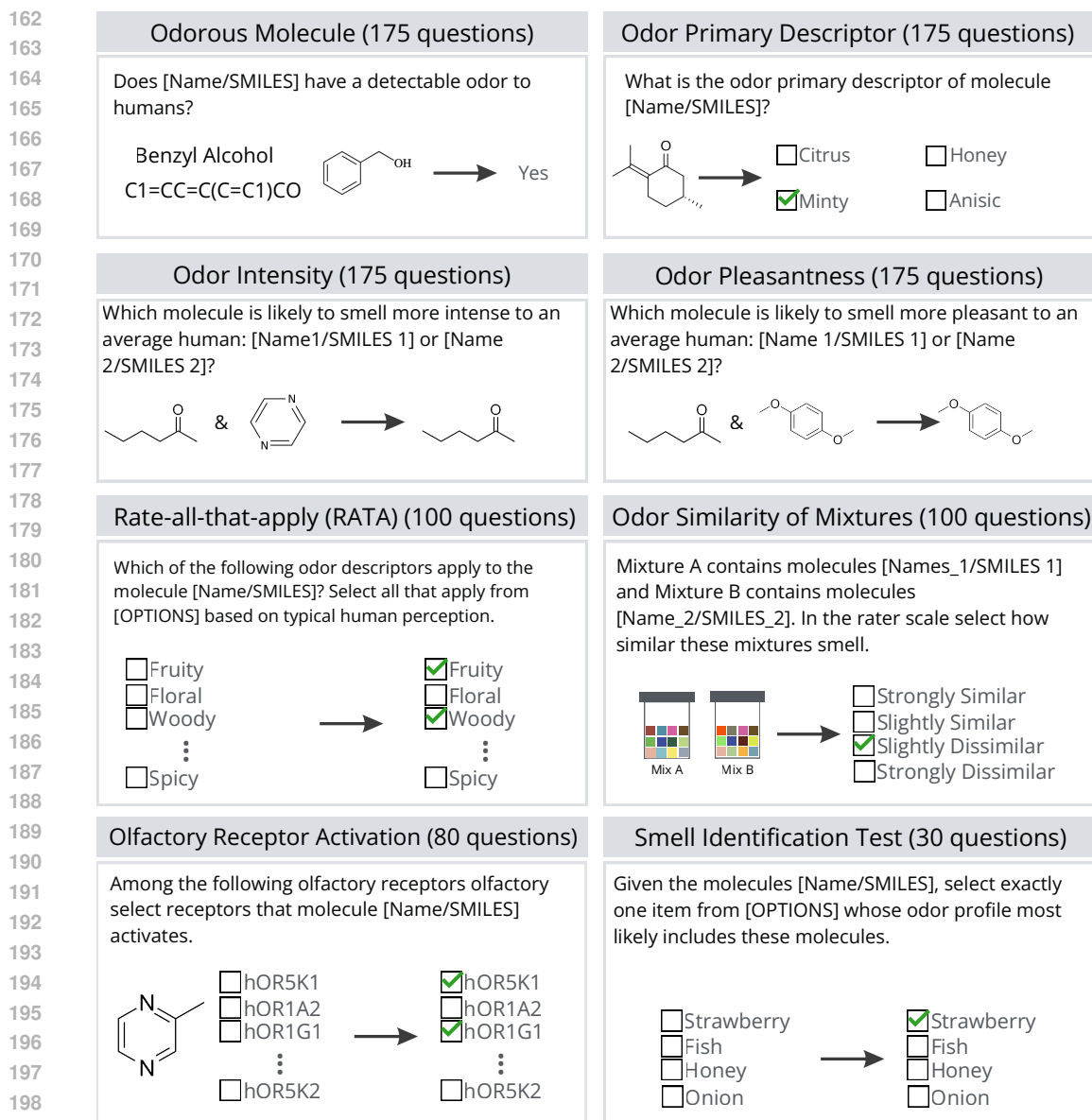


Figure 1: Olfactory Perception (OP) benchmark task description and output format across tasks.

four bins of a categorical rating scale, including: strongly similar, slightly similar, slightly dissimilar, and strongly dissimilar. This task deals with mixture-level olfactory perception, and we classify it as intermediate-difficulty.

Olfactory Receptor Activation (ORA). In this task, we aim to assess model capability in identifying olfactory receptor (OR) activation for pure molecules. We use a multilabel task (80 questions) that asks which receptors (from a candidate set of 4–10 human OR gene IDs) are activated by a given molecule. The molecules-OR pairs were obtained from the M2OR dataset Lalis et al. (2024). We constrained our selection of molecules-OR pairs to cases where 4 to 10 ORs activation were experimentally studied, and 1 to 10 ORs were observed to be activated. We consider ORA a hard task.

Smell Identification Test (SIT). We curated a multiple-choice smell identification task (30 questions) with molecular mixture information. The model selects the most likely odor source from a four-option

216 set given a molecular profile. We obtained the list of volatile organic compounds for different food
217 items from the Leibniz-Isbtum odorant database Kreissl et al. (2022). The items studied include the
218 odor profile for mango, peanut, hazelnut, tomato, apple, walnut, raspberry, peach, honey, parsley,
219 grapefruit, pineapple, strawberry, apricot, rice, grape, popcorn, orange, cheese, melon, leather,
220 chocolate, coffee, onion, fish, beer, whisky, red wine, prawn, and bread. We consider the SIT a hard
221 task.

222 **Multilingual Translation.** We generated multilingual versions of the RATA tasks by translating
223 the English odor-descriptor vocabulary and the UI/instruction text into each target language using
224 GPT-5 . 2. We prompted the model to translate in a perfumery/sensory context and to produce brief,
225 natural descriptors; we used high-reasoning settings for disambiguation. Descriptor translations were
226 normalized to lowercase for consistency, while UI strings were translated separately to preserve
227 natural phrasing. To reduce hallucinated or awkward terms, we applied a lightweight two-step quality
228 control check for any uncertain cases: the model proposes candidate single-word descriptors and then
229 selects the best option based on whether it is a real word in the target language and whether it fits
230 olfactory usage. If no suitable translation is found, we fall back to the closest candidate; if generation
231 fails entirely, we retain the original English descriptor. For languages where compounds are standard
232 (e.g., German, Swedish), compound forms were allowed as long as they remained whitespace-free.

233 4 EXPERIMENTS

234 To comprehensively evaluate olfactory reasoning capabilities across the current LLM landscape, we
235 select models spanning multiple providers, architectures, and reasoning configurations. A key dimen-
236 sion of our evaluation is the effect of reasoning budget on olfactory task performance, motivated by
237 findings in chemistry benchmarks showing that extended reasoning substantially improves molecular
238 understanding Runcie et al. (2025). We organize our evaluation into two groups:

- 241 • **Closed-Source Models:** We evaluate models from five providers. From OpenAI Achiam et al.
242 (2023), we test the reasoning models o3 (high) and o4-mini (high), the GPT-5 family Singh et al.
243 (2025) at two reasoning levels (low, high), GPT-5 Pro, GPT-5.2 Pro, and GPT-OSS 120B Agarwal
244 et al. (2025) at high reasoning settings. From Google, we evaluate Gemini 2.5 Pro Comanici et al.
245 (2025) at three reasoning budgets (8K, 16K, and 32K tokens) to isolate the effect of reasoning
246 depth. From xAI, we include Grok 3 Mini at low and high reasoning and Grok 4.1 Fast. From
247 Anthropic, we evaluate Claude Sonnet 4.5, Claude Opus 4.5 Anthropic (2025), and Claude Opus
248 4.6 Anthropic (2026) at two reasoning budgets (high, max).
- 249 • **Open-Source Models:** We evaluate DeepSeek Reasoner Guo et al. (2025) at three reasoning
250 budgets (8K, 16K, and 32K tokens) and Llama-3.3-70B-Instruct Dubey et al. (2024). These models
251 provide a baseline for open-weight performance on olfactory tasks.

252 In total, we evaluate 21 model configurations across 6 providers. All models were queried through
253 their respective provider APIs without enabling web search, tool use, or any external retrieval capabil-
254 ities; by default, none of the APIs grant models access to external resources during inference. All
255 models are prompted identically using both isomeric SMILES and compound name representations;
256 full prompt templates and worked examples are provided in Appendix C. For reasoning models,
257 we systematically vary the reasoning budget to analyze the relationship between computational
258 deliberation and olfactory task accuracy.

259 **Evaluation Metrics.** We adopt task-appropriate metrics to capture performance across the diverse
260 question formats in the OP benchmark. For single-answer tasks (OC, OPD, OIn, OPI, OS, SIT),
261 we report *any-overlap accuracy*: a question is scored correct if the set of tokens extracted from
262 the model’s response intersects the ground-truth token set. All prompts explicitly instruct models
263 to respond with only the selected answer and no additional commentary; a post-hoc cleaning step
264 verified compliance and corrected the very few cases (<1% of responses, limited to Claude 4.6 and
265 Grok Mini reasoning traces) where extraneous tokens appeared. For multi-answer tasks (RATA,
266 ORA), we employ *per-question multilabel F1*, which provides partial credit when models correctly
267 identify a subset of applicable labels and penalises both spurious and missing predictions. For the
268 three continuous-rating tasks (OIn, OPI, OS), we additionally compute Pearson correlations between
269 predicted and human psychophysical ratings. Overall accuracy is the unweighted arithmetic mean of
the eight per-task scores. More details on the answer extraction pipeline are provided in Appendix D.

Table 1: Olfactory Perception (OP) benchmark results using compound name prompts. Accuracy (%) is reported per task and overall (unweighted mean). **OC** = Odor Classification ($n=175$), **OPD** = Odor Primary Descriptor ($n=175$), **OIn** = Odor Intensity ($n=175$), **OPI** = Odor Pleasantness ($n=175$), **RATA** = Rate-All-That-Apply ($n=100$), **OS** = Odor Similarity ($n=100$), **ORA** = Olfactory Receptor Activation ($n=80$), **SIT** = Smell Identification Test ($n=30$). **Bold** = best overall; underlined = second best.

Model	Reasoning	Simple ($N=700$)				Intermediate ($N=200$)		Hard ($N=110$)		Overall
		OC	OPD	OIn	OPI	RATA	OS	ORA	SIT	
Closed-Source										
GPT-5	high	89.7	73.7	66.3	71.4	36.4	<u>34.0</u>	40.8	<u>76.7</u>	61.1
GPT-5	low	90.3	72.0	70.9	71.4	30.8	28.0	40.4	<u>73.3</u>	59.6
GPT-5 Pro	high	92.0	73.1	71.4	70.9	36.1	29.0	42.5	80.0	61.9
GPT-5.2 Pro	high	88.6	76.0	<u>72.6</u>	72.0	34.6	28.0	52.8	73.3	62.2
GPT-OSS-120B	high	82.9	60.6	<u>65.1</u>	72.0	25.1	<u>34.0</u>	35.6	56.7	54.0
o3	high	89.1	70.9	68.0	70.3	31.5	32.0	42.4	70.0	59.3
o4-mini	high	88.6	65.7	69.1	73.7	29.0	32.0	40.5	73.3	59.0
Gemini 2.5 Pro	16K	89.7	<u>78.9</u>	68.0	72.6	30.3	31.0	39.7	66.7	59.6
Gemini 2.5 Pro	32K	87.4	78.3	66.9	72.0	34.0	27.0	42.4	70.0	59.7
Gemini 2.5 Pro	8K	88.6	80.0	65.7	73.7	31.5	29.0	37.9	63.3	58.7
Grok 3 Mini	high	81.7	72.0	68.0	73.7	37.0	22.0	41.9	66.7	57.9
Grok 3 Mini	low	81.7	73.1	66.3	72.6	36.0	18.0	37.2	73.3	57.3
Grok 4.1 Fast	default	88.6	67.4	66.9	70.3	35.5	33.0	31.1	73.3	58.3
Claude Opus 4.5	high	92.0	76.6	71.4	73.1	42.2	25.0	45.8	70.0	62.0
Claude Opus 4.6	high	<u>91.4</u>	78.3	71.4	74.9	<u>40.0</u>	26.0	49.6	73.3	<u>63.1</u>
Claude Opus 4.6	max	92.0	77.7	74.9	<u>74.3</u>	38.9	26.0	<u>51.1</u>	80.0	64.4
Claude Sonnet 4.5	—	89.1	67.4	66.9	71.4	34.9	29.0	38.4	80.0	59.6
Open-Source										
DeepSeek Reasoner	16K	79.4	69.7	68.6	74.9	36.0	25.0	29.1	70.0	56.6
DeepSeek Reasoner	32K	81.1	73.7	69.7	73.1	33.1	32.0	31.6	73.3	58.5
DeepSeek Reasoner	8K	80.6	70.3	71.4	72.6	34.7	35.0	30.5	63.3	57.3
Llama 3.3 70B	—	83.4	60.6	68.0	72.0	26.8	29.0	35.0	46.7	52.7

5 RESULTS

In this section, we evaluate a diverse set of state-of-the-art LLMs on the OP benchmark to quantify current olfactory reasoning ability and identify systematic failure modes. We report results across all task categories and compare prompting conditions (compound names vs. isomeric SMILES) to isolate the effect of molecular representation on performance.

Overall Performance Table 1 and Figure 2 present benchmark performance across all evaluated models. The best-performing configuration, **Claude Opus 4.6 (max)**, attains 64.4% overall accuracy with compound name prompts, followed by **Claude Opus 4.6 (high)** at 63.1%, **GPT-5.2 Pro** at 62.2%, and **Claude Opus 4.5** at 62.0%. These scores substantially exceed chance levels for each task category (Figure 2b, dashed lines), yet remain far from ceiling performance, indicating that, while frontier models encode meaningful olfactory knowledge, substantial room for improvement remains.

A clear performance hierarchy emerges across providers. Anthropic and OpenAI frontier models occupy the top positions, with all configurations except GPT-OSS-120B surpassing 59% accuracy. The progression within the Claude family—from Sonnet 4.5 (59.6%) through Opus 4.5 (62.0%) to Opus 4.6 (max) (64.4%)—reflects consistent improvements associated with increased model capability and reasoning depth. OpenAI’s GPT-5 variants demonstrate competitive performance, with GPT-5.2 Pro trailing Claude Opus 4.6 (max) by 2.2 percentage points. The o-series reasoning models perform well, with o3 (high) at 59.3% and o4-mini (high) at 59.0%. Google’s Gemini 2.5 Pro (58.7–59.7%) and xAI’s Grok family (57.3–58.3%) occupy the mid-tier, while DeepSeek Reasoner reaches 56.6–58.5%. The open-source Llama 3.3 70B lags behind all proprietary alternatives at

52.7%, underscoring a persistent capability gap between open-weight and closed-source systems for specialized perceptual reasoning.

Performance varies considerably across task categories. Simple tasks (OC, OPD, OIn, OPI) produce the highest accuracy, with top models reaching 92.0% on odor classification and 80.0% on primary descriptor identification. Intermediate tasks (RATA, OS) prove more demanding, with best performance limited to 42.2% and 35.0% respectively. Hard tasks exhibit the widest variance: SIT reaches 80.0% for three models, benefiting from knowledge about foods and beverages, whereas ORA peaks at 52.8% (GPT-5.2 Pro), still a challenging task due to the specialized biochemistry knowledge required. A detailed question-level difficulty analysis is provided in Appendix E.1 (Figure 5) for single-label tasks and Appendix E.2 (Figure 6) for the multi-label tasks RATA and ORA.

Isomeric SMILES vs. Compound Name Prompts A consistent and substantial performance gap separates the two molecular representation formats (Figure 2a). Across all 21 model configurations, compound name prompts outperform isomeric SMILES notation, with improvements ranging from

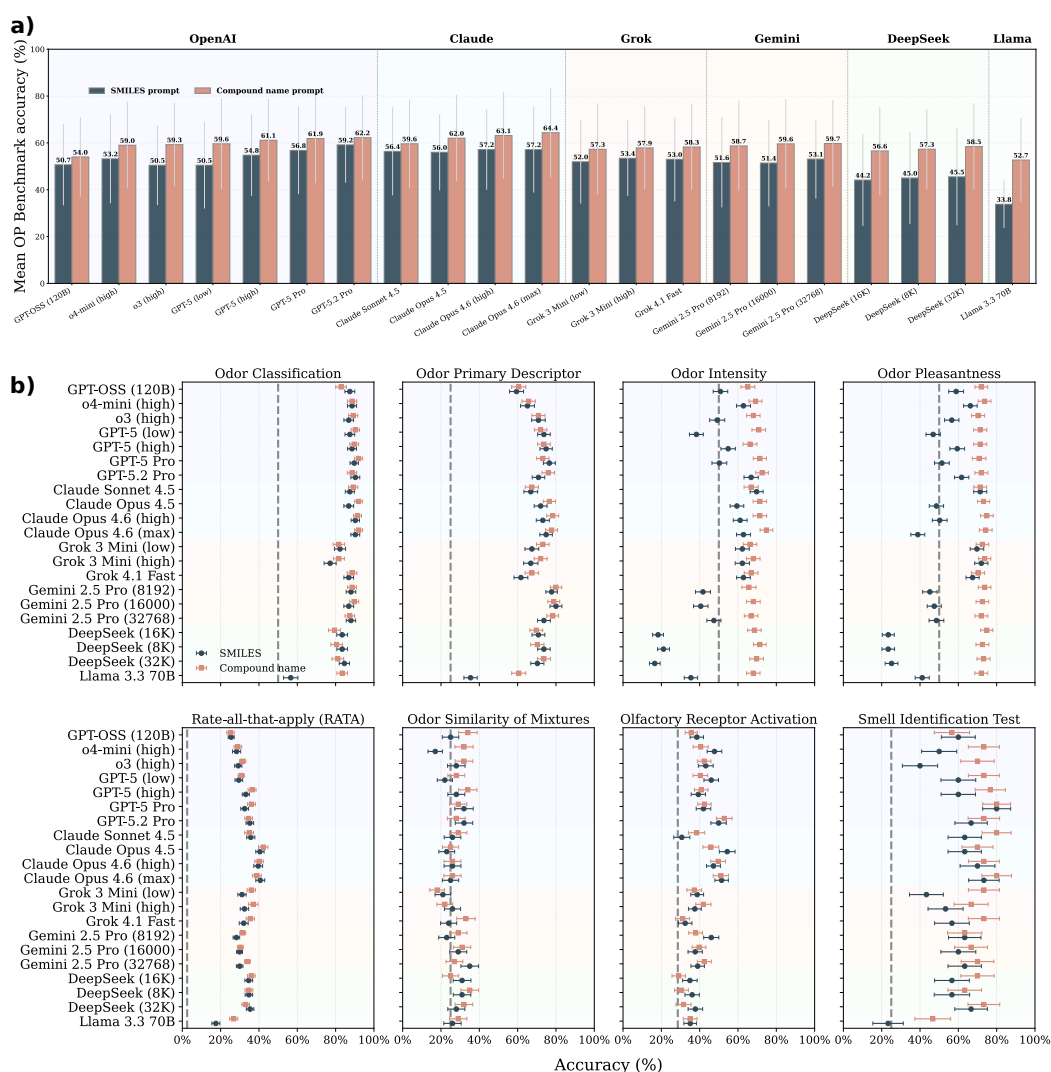


Figure 2: **Olfactory Perception (OP) benchmark performance.** (a) Overall mean accuracy per model for isomeric SMILES and compound name prompts. Error bars: 95% confidence intervals. (b) Per-category accuracy for each model under both prompts. Horizontal error bars are bootstrap standard deviations of the mean accuracy. Gray dashed lines indicate chance baselines.

2.4 percentage points (Claude Opus 4.5: 59.6% \rightarrow 62.0%) to 18.9 percentage points (Llama 3.3 70B: 33.8% \rightarrow 52.7%). DeepSeek Reasoner (8K) shows a 12.3-point improvement (45.0% \rightarrow 57.3%). The mean improvement across models is approximately 7 percentage points, suggesting that current LLMs access olfactory knowledge primarily through lexical associations rather than structural molecular reasoning. Frontier reasoning models display the narrowest gaps, with GPT-5 Pro, GPT-5.2 Pro, and Claude Opus 4.5 each retaining over 95% of their compound name accuracy under SMILES prompts. Conversely, Llama 3.3 70B exhibits the largest disparity, with isomeric SMILES accuracy barely exceeding chance. Per-task analysis (Figure 2b) reveals that the isomeric SMILES/name gap varies substantially by task category. OC exhibits the smallest gaps, potentially because odorousness correlates with molecular properties (volatility, molecular weight) that can be partially inferred from isomeric SMILES structure. OIn, OPI, and SIT show the largest gaps, consistent with tasks where identifying the target molecule or food source is critical and far easier from a name than from structural notation. In contrast, multi-label tasks (RATA, ORA) show modest gaps, suggesting that descriptor-level knowledge is similarly accessible, or similarly limited, under both representations.

Effect of Reasoning Budget Systematic variation of reasoning token budgets reveals consistent but modest performance gains. Claude Opus 4.6 shows clear scaling from high (63.1%) to max (64.4%), and the GPT-5 family improves from 59.6% (low) to 62.2% (GPT-5.2 Pro). Similar patterns appear across other providers, though DeepSeek Reasoner shows a non-monotonic pattern, suggesting optimal reasoning budgets may be task-dependent. Overall, no model gains more than approximately 2 percentage points from extended reasoning, contrasting with chemistry benchmarks where reasoning provided substantially larger gains Runcie et al. (2025).

Fine-Grained Performance Analysis Beyond categorical accuracy, we assess whether models capture the continuous structure of human olfactory perception. Figure 3a presents Pearson correlations between model-predicted ratings and human psychophysical measurements: the best models reach $r \approx 0.55$ for OIn, approaching specialized model performance Keller et al. (2017); Satarifard et al. (2025) (red dashed line); OPI correlations are higher ($r \approx 0.60$); and OS exhibits the weakest correlations ($r \approx 0.35$), confirming that models struggle to integrate perceptual information across molecules. Isomeric SMILES and compound name prompts produce similar correlations across all three dimensions, with a slightly larger gap for OIn. Turning to the multi-label tasks, RATA and ORA exhibit distinct difficulty profiles: RATA F1 scores are mainly bell-shaped, indicating partial credit on most questions, while ORA is bimodal: models either possess the relevant receptor-ligand knowledge or lack it entirely (Appendix E.2, Figure 6). Per-label analysis (Figure 3b) reveals a clear divide for RATA: descriptors with well-established functional-group associations (e.g., sulfurous, floral, fruity Genva et al. (2019)) achieve the highest mean F1, while descriptors such as spicy, fresh, and tropical prove nearly impossible. A detailed case study of this failure mode is provided in Appendix E.3. For ORA, per-label difficulty varies sharply across receptors (Figure 3b): the wildtype hOR2W1 (appearing in 30 label assignments) is reliably predicted at F1 above 0.6, and the M81V variant is moderately well-predicted, but hOR2W1_D296N (24) and hOR52D1 (5) sit near zero for most models. The hOR52D1 gap likely reflects data scarcity, while D296N reveals a specific knowledge error: Claude Opus 4.6 (max) never predicts this receptor across all 24 label appearances, whereas GPT-5.2 Pro correctly identifies it in 19 cases. Inter-model variance is substantially higher for receptor labels than for semantic descriptors, indicating that receptor knowledge is more idiosyncratic across model families; a detailed analysis of these knowledge gaps is provided in Appendix F.2 (Figure 10).

Systematic Failure Modes Two qualitatively distinct failure mechanisms underlie the hardest categories. For OS, models use molecular overlap as a proxy for perceptual similarity: accuracy reaches $\sim 85\%$ when similar mixtures share many molecules but drops to near 0% when similar mixtures share few, and all models exhibit a systematic bias toward predicting dissimilarity; Claude models assign ‘‘Slightly Dissimilar’’ to the vast majority of mixtures, even though perceptually similar mixtures often share zero molecules. This heuristic renders mixture-level olfactory similarity fundamentally beyond current LLM capabilities; a detailed analysis is provided in Appendix F.1 (Figures 8, 7, 9). Additionally, Claude Opus 4.6’s safety filter refuses OC questions about hazardous compounds (e.g., nerve agents), illustrating a tension between safety alignment and scientific evaluation (Appendix F.3).

Multilingual Evaluation To probe whether olfactory knowledge is language-specific, we translated the RATA task prompts into 21 languages spanning over six language families and evaluated seven models per language. Figure 4a shows that English achieves the highest mean per-question F1, followed closely by French, Spanish, and Russian, while non-Indo-European languages (Korean,

Chinese, Swahili) cluster at the lower end. Figure 4b presents language-family AUROC curves for DeepSeek (32K): Germanic and Romance languages achieve the highest family-level AUROC (0.778 and 0.774, respectively), while the remaining families score lower (0.725–0.749). An ensemble aggregating votes across all 21 languages achieves AUROC = 0.828, suggesting that multilingual aggregation captures complementary olfactory knowledge. Figure 4c shows per-model AUROC using an all-language vote-fraction ensemble: Gemini 2.5 Pro (32K) leads at 0.864, and a cross-model

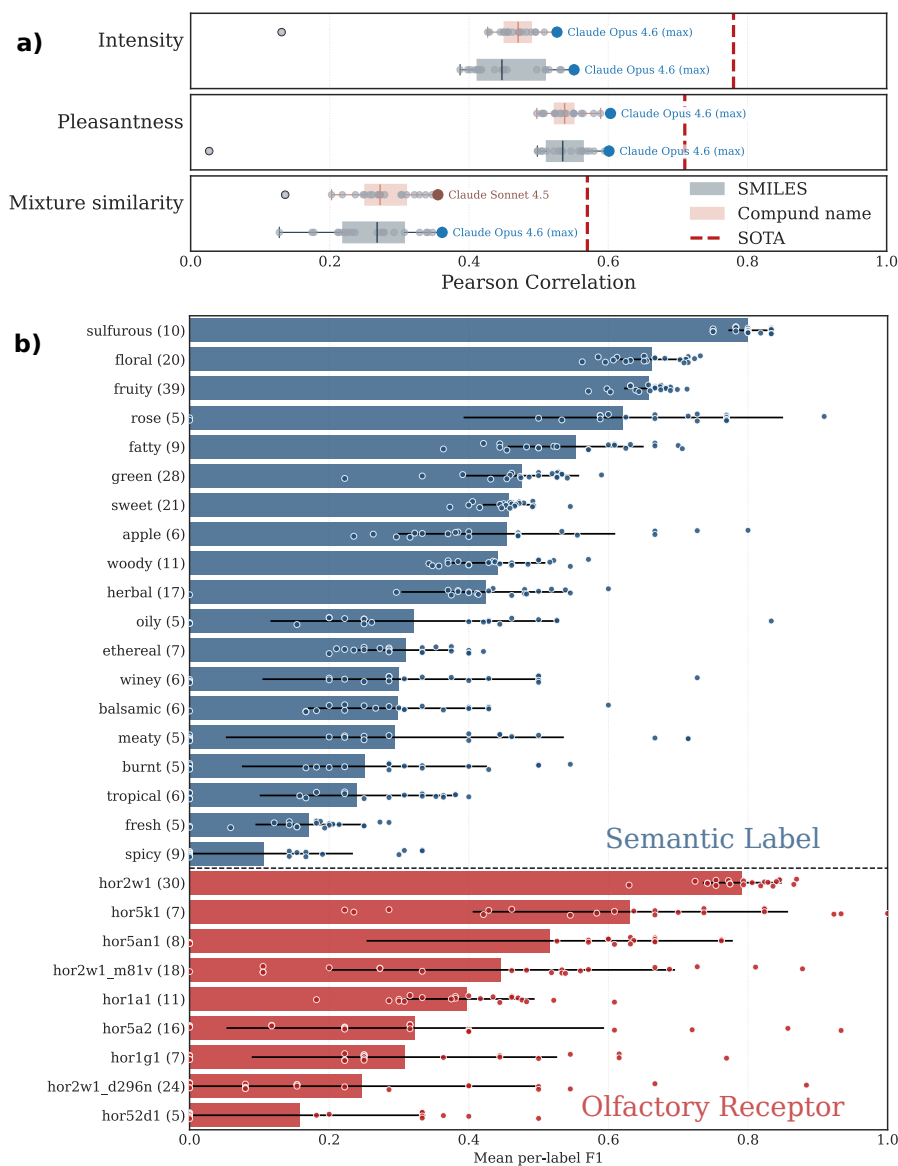


Figure 3: **Predictions correlations and performance difficulty across olfactory tasks.** (a) Pearson correlations across models for three continuous-rating categories (odor intensity, odor pleasantness, and mixture similarity). Each row shows the distribution of model correlations for isomeric SMILES and Compound name prompts; gray points represent individual models, and the best-performing model per prompt/category is highlighted with a colored circle. State-of-the-art reference performance is indicated by red dashed lines. (b) Combined label-difficulty ranking for RATA (blue) and olfactory-receptor activation (red). For each label, per-model F1 scores are computed in a multilabel setting; the bar shows the mean F1 across models with an error bar (standard deviation), while points represent individual model values.

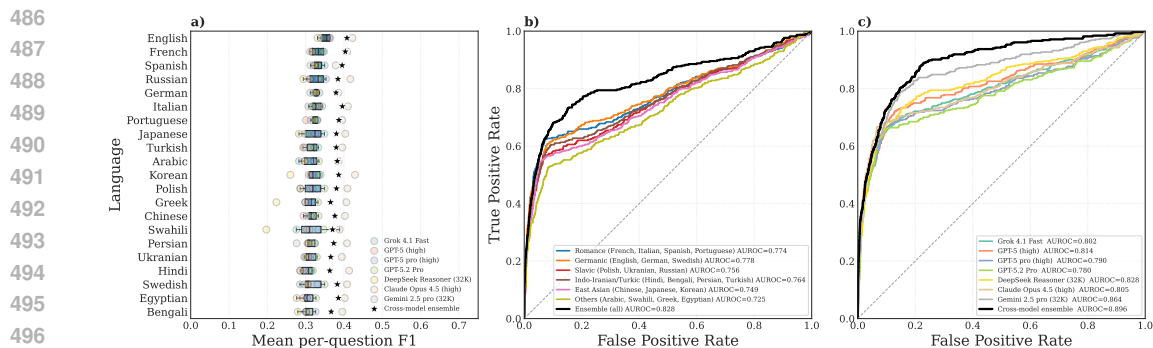


Figure 4: **Multilingual RATA performance for compound name prompt.** (a) Per-language distributions of model performance, summarized as mean per-question multilabel F1, colored points denote individual models and the black star denotes a cross-model ensemble. (b) DeepSeek (32K) AUROC by language family using per-label vote fractions across languages within each family (East Asian and Others are not same language families), black curve represents an ensemble of all languages. (c) AUROC for each model using an all-language majority vote ensemble, plus a cross-model ensemble pooling all models from all languages. Dashed diagonal indicates chance.

ensemble combining all seven models reaches 0.896, indicating that model diversity further improves multilingual olfactory prediction beyond language diversity alone.

6 CONCLUSION, LIMITATIONS, AND FUTURE WORK

Here, we introduce the Olfactory Perception (OP) benchmark, a structured evaluation suite for testing whether large language models can reason about smell from molecular information and real-world odor sources. Across a broad set of olfactory tasks, we observe emerging but incomplete capabilities: the best-performing systems achieve moderate accuracy, yet performance varies substantially by task type. A consistent pattern is the gap between prompts using compound names versus isomeric SMILES, suggesting current LLMs often succeed via lexical associations rather than robust structural reasoning. Taken together, our results position olfaction as a challenging and underexplored modality for LLM evaluation, and provide a concrete benchmark for measuring progress toward models that can connect chemistry to sensory perception. The benchmark makes deliberate design trade-offs: we use standardized descriptor vocabularies and discrete response formats, which enable reproducible automatic evaluation but do not capture the full richness of odor perception (e.g., contextual effects, individual and cultural variation, nuanced free-form descriptions). As with other instruction-following benchmarks, results can be influenced by prompting and output formatting sensitivity, despite using consistent templates. Multilingual variants rely on automated translation with validation, which may introduce subtle connotation shifts. Benchmark accuracy should be interpreted as agreement with established labels rather than as evidence of mechanistic olfactory understanding. Several directions could strengthen future work. On the data side, expanding to mixtures, concentration effects, temporal dynamics, broader receptor-level coverage, and multimodal inputs (e.g., verbal descriptions, facial expressions during smelling) would better reflect real olfaction. Incorporating human evaluation and cross-cultural annotation could capture partial correctness and reduce brittleness to surface forms. On the modeling side, hybrid systems combining LLMs with cheminformatics tools, methods encouraging structure-based reasoning (e.g., adversarial splits, representation-invariant prompting), and low-cost strategies such as prompt repetition (Leviathan et al. (2025)) are promising directions. Finally, extending to generative settings, such as proposing molecules with target odor profiles under safety and synthesizability constraints, would connect benchmark performance to applications in fragrance, flavor, and environmental monitoring.

REFERENCES

- 540
541
542 Erratum for the report humans can discriminate more than 1 trillion olfactory stimuli by c. bushdid et
543 al. *Science*, 383(6685):eado6457, 2024. doi: 10.1126/science.ado6457. URL [https://www.
544 science.org/doi/abs/10.1126/science.ado6457](https://www.science.org/doi/abs/10.1126/science.ado6457).
- 545 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
546 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
547 *arXiv preprint arXiv:2303.08774*, 2023.
- 548 Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K
549 Arora, Yu Bai, Bowen Baker, Haiming Bao, et al. gpt-oss-120b & gpt-oss-20b model card. *arXiv
550 preprint arXiv:2508.10925*, 2025.
- 551 Anthropic. Claude Opus 4.5 system card. Technical report, 2025. URL [https://www.
552 anthropic.com/claude-opus-4-5-system-card](https://www.anthropic.com/claude-opus-4-5-system-card).
- 553 Anthropic. Claude Opus 4.6 system card. Technical report, 2026. URL [https://www-cdn.
554 anthropic.com/0dd865075ad3132672ee0ab40b05a53f14cf5288.pdf](https://www-cdn.anthropic.com/0dd865075ad3132672ee0ab40b05a53f14cf5288.pdf).
- 555 Jennifer Bartsch, Erik Uhde, and Tunga Salthammer. Analysis of odour compounds
556 from scented consumer products using gas chromatography-mass spectrometry and gas
557 chromatography-olfactometry. *Analytica chimica acta*, 904:98–106, 2016. URL [https:
558 //api.semanticscholar.org/CorpusID:7550203](https://api.semanticscholar.org/CorpusID:7550203).
- 559 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
560 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
561 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler,
562 Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott
563 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya
564 Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL [https:
565 //arxiv.org/abs/2005.14165](https://arxiv.org/abs/2005.14165).
- 566 Linda B. Buck and Richard Axel. A novel multigene family may encode odorant receptors:
567 a molecular basis for odor recognition. *Cell*, 65 1:175–87, 1991. URL [https://api.
568 semanticscholar.org/CorpusID:6548928](https://api.semanticscholar.org/CorpusID:6548928).
- 569 Caroline Bushdid, Marcelo O Magnasco, Leslie B Vosshall, and Andreas Keller. Humans can
570 discriminate more than 1 trillion olfactory stimuli. *Science*, 343(6177):1370–1372, 2014.
- 571 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Ka-
572 plan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen
573 Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray,
574 Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens
575 Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis,
576 Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas
577 Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher
578 Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford,
579 Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario
580 Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language
581 models trained on code, 2021. URL <https://arxiv.org/abs/2107.03374>.
- 582 Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit
583 Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier
584 with advanced reasoning, multimodality, long context, and next generation agentic capabilities.
585 *arXiv preprint arXiv:2507.06261*, 2025.
- 586 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
587 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
588 *arXiv e-prints*, pp. arXiv–2407, 2024.
- 589 Sergio Esteban-Romero, Iván Martín-Fernández, Manuel Gil-Martín, and Fernando Fernández-
590 Martínez. Synthesizing olfactory understanding: Multimodal language models for image–text
591 smell matching. *Symmetry*, 17(8):1349, 2025.

- 594 Makoto Fukushima, Shusuke Eshita, and Hiroshige Fukuhara. Advancements and limitations of llms
595 in replicating human color-word associations. *Discover Artificial Intelligence*, 5(1):64, 2025.
596
- 597 Xinyu Ge, Yongjie Zhou, Qing Li, Yuqing Tan, Yongkang Luo, and Hui Hong. Machine learning for
598 food flavor prediction and regulation: models, data integration, and future perspectives. *Journal*
599 *of advanced research*, 2025. URL [https://api.semanticscholar.org/CorpusID:
600 282185617](https://api.semanticscholar.org/CorpusID:282185617).
- 601 Manon Genva, Thierry Kenne Kemene, Magali Deleu, Laurence Lins, and Marie-Laure Fauconnier. Is
602 it possible to predict the odor of a molecule on the basis of its structure? *International journal of*
603 *molecular sciences*, 20(12):3018, 2019.
- 604 Richard C Gerkin and Jason B Castro. The number of olfactory stimuli that humans can discriminate
605 is still unknown. *eLife*, 4:e08127, jul 2015. ISSN 2050-084X. doi: 10.7554/eLife.08127. URL
606 <https://doi.org/10.7554/eLife.08127>.
607
- 608 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
609 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
610 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 611 Taicheng Guo, Bozhao Nan, Zhenwen Liang, Zhichun Guo, Nitesh Chawla, Olaf Wiest, Xiangliang
612 Zhang, et al. What can large language models do in chemistry? a comprehensive benchmark on
613 eight tasks. *Advances in Neural Information Processing Systems*, 36:59662–59688, 2023.
- 614 Yang Han, Ziping Wan, Lu Chen, Kai Yu, and Xin Chen. From generalist to specialist: A survey of
615 large language models for chemistry, 2024. URL <https://arxiv.org/abs/2412.19994>.
616
- 617 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob
618 Steinhardt. Measuring massive multitask language understanding. In *International Confer-*
619 *ence on Learning Representations*, 2021. URL [https://openreview.net/forum?id=
620 d7KBjmI3GmQ](https://openreview.net/forum?id=d7KBjmI3GmQ).
- 621 International Fragrance Association. The ifra fragrance ingredient glossary
622 (fig), 2020. URL [https://d3t14p1xronwr0.cloudfront.net/docs/
623 ifra-fragrance-ingredient-glossary-april-2020.pdf](https://d3t14p1xronwr0.cloudfront.net/docs/ifra-fragrance-ingredient-glossary-april-2020.pdf).
624
- 625 Andreas Keller, Richard C Gerkin, Yuanfang Guan, Amit Dhurandhar, Gabor Turu, Bence Szalai,
626 Joel D Mainland, Yusuke Ihara, Chung Wen Yu, Russ Wolfinger, et al. Predicting human olfactory
627 perception from chemical features of odor molecules. *Science*, 355(6327):820–826, 2017.
- 628 J. Kreissl, V. Mall, P. Steinhaus, and M. Steinhaus. Leibniz-lsb@tum odorant
629 database, version 1.2, 2022. URL [https://www.leibniz-lsb.de/en/databases/
630 leibniz-lsbtum-odorant-database](https://www.leibniz-lsb.de/en/databases/leibniz-lsbtum-odorant-database).
631
- 632 Murathan Kurfalı, Jonas K Olofsson, and Thomas Hörberg. Enhancing multimodal language models
633 with olfactory information. 2023.
- 634 Murathan Kurfalı, Pawel Herman, Stephen Pierzchajlo, Jonas Olofsson, and Thomas Hörberg.
635 Representations of smells: The next frontier for language models? *Cognition*, 264:106243, 2025.
636
- 637 Maxence Lalis, Matej Hladiš, Samar Abi Khalil, Loïc Briand, Sébastien Fiorucci, and Jérémie Topin.
638 M2or: a database of olfactory receptor–odorant pairs for understanding the molecular mechanisms
639 of olfaction. *Nucleic Acids Research*, 52(D1):D1370–D1379, 2024.
- 640 Brian K Lee, Emily J Mayhew, Benjamin Sanchez-Lengeling, Jennifer N Wei, Wesley W Qian,
641 Kelsie A Little, Matthew Andres, Britney B Nguyen, Theresa Moloy, Jacob Yasonik, et al. A
642 principal odor map unifies diverse tasks in olfactory perception. *Science*, 381(6661):999–1006,
643 2023.
- 644 Yaniv Leviathan, Matan Kalman, and Yossi Matias. Prompt repetition improves non-reasoning llms.
645 *arXiv preprint arXiv:2512.14982*, 2025.
646
- 647 Joel Mainland, Yun Li, Ting Zhou, Wen Liu, and Hiroaki Matsunami. Human olfactory receptor
responses to odorants. *Scientific Data*, 2:150002, 02 2015. doi: 10.1038/sdata.2015.2.

- 648 Asifa Majid. Human olfaction at the intersection of language, culture, and biology. *Trends in*
649 *Cognitive Sciences*, 25(2):111–123, 2021. ISSN 1364-6613. doi: [https://doi.org/10.1016/j.tics.](https://doi.org/10.1016/j.tics.2020.11.005)
650 2020.11.005. URL [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S1364661320302771)
651 [S1364661320302771](https://www.sciencedirect.com/science/article/pii/S1364661320302771).
- 652 Asifa Majid and Niclas Burenhult. Odors are expressible in language, as long as you speak the
653 right language. *Cognition*, 130(2):266–270, 2014. ISSN 0010-0277. doi: [https://doi.org/](https://doi.org/10.1016/j.cognition.2013.11.004)
654 [10.1016/j.cognition.2013.11.004](https://doi.org/10.1016/j.cognition.2013.11.004). URL [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/pii/S001002771300214X)
655 [article/pii/S001002771300214X](https://www.sciencedirect.com/science/article/pii/S001002771300214X).
- 656 Asifa Majid, Niclas Burenhult, Marcus Stensmyr, Josje Valk, and Bill Hansson. Olfactory language
657 and abstraction across cultures. *Philosophical Transactions of the Royal Society B: Biological*
658 *Sciences*, 373, 06 2018. doi: 10.1098/rstb.2017.0139.
- 659 Haitao Mao, Linlin Liu, Jian Du, and Rafiqul Gani. A machine learning based computer-aided molec-
660 ular design/screening methodology for fragrance molecules. *Computers & Chemical Engineering*,
661 115, 04 2018. doi: 10.1016/j.compchemeng.2018.04.018.
- 662 Raja Marjieh, Iliia Sucholutsky, Pol van Rijn, Nori Jacoby, and Thomas L Griffiths. Large language
663 models predict human sensory judgments across six modalities. *Scientific Reports*, 14(1):21445,
664 2024.
- 665 Emily J Mayhew, Charles J Arayata, Richard C Gerkin, Brian K Lee, Jonathan M Magill, Lindsey L
666 Snyder, Kelsie A Little, Chung Wen Yu, and Joel D Mainland. Transport features predict if a
667 molecule is odorous. *Proceedings of the National Academy of Sciences*, 119(15):e2116576119,
668 2022.
- 669 Adrian Mirza, Nawaf Alampara, Sreekanth Kunchapu, Martiño Ríos-García, Benedict Emoekabu,
670 Aswanth Krishnan, Tanya Gupta, Mara Schilling-Wilhelmi, Macjonathan Okereke, Anagha Aneesh,
671 et al. A framework for evaluating the chemical knowledge and reasoning abilities of large language
672 models against the expertise of chemists. *Nature Chemistry*, pp. 1–8, 2025.
- 673 Andrew Oxenham, Christophe Micheyl, Michael Keebler, Adam Loper, and Sébastien Santurette.
674 Pitch perception beyond the traditional existence region of pitch. *Proceedings of the National*
675 *Academy of Sciences of the United States of America*, 108:7629–34, 05 2011. doi: 10.1073/pnas.
676 1015291108.
- 677 Aharon Ravia, Kobi Snitz, Danielle Honigstein, Maya Finkel, Rotem Zirler, Ofer Perl, Lavi Secundo,
678 Christophe Laudamiel, David Harel, and Noam Sobel. A measure of smell enables the creation of
679 olfactory metamers. *Nature*, 588(7836):118–123, 2020.
- 680 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani,
681 Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In
682 *First Conference on Language Modeling*, 2024.
- 683 Nicholas T Runcie, Charlotte M Deane, and Fergus Imrie. Assessing the chemical intelligence of
684 large language models. *Journal of Chemical Information and Modeling*, 2025.
- 685 Pinaki Saha, Mrityunjay Sharma, Sarabeshwar Balaji, Aryan Amit Barsainyan, Ritesh Kumar,
686 Volker Steuber, and Michael Schmucker. Dense sense: a novel approach utilizing electron density
687 augmented machine learning paradigm to understand the complex odour landscape. *Digital*
688 *Discovery*, 4(11):3339–3350, 2025.
- 689 Vahid Satarifard, Laura Sisson, Yikun Han, Pedro Ilídio, Matej Hladiš, Maxence Lalis, Xuebo Song,
690 Wenjie Yin, Aharon Ravia, CiCi Xingyu Zheng, et al. High-fidelity tuning of olfactory mixture
691 distances in the perceptual space of smell through a community effort. *bioRxiv*, pp. 2025–12, 2025.
- 692 Lavi Secundo, Kobi Snitz, and Noam Sobel. The perceptual logic of smell. *Current opinion in*
693 *neurobiology*, 25C:107–115, 01 2014. doi: 10.1016/j.conb.2013.12.010.
- 694 Aaditya Singh, Adam Fry, Adam Perelman, Adam Tart, Adi Ganesh, Ahmed El-Kishky, Aidan
695 McLaughlin, Aiden Low, AJ Ostrow, Akhila Ananthram, et al. Openai gpt-5 system card. *arXiv*
696 *preprint arXiv:2601.03267*, 2025.

- 702 Kobi Snitz, Adi Yablonka, Tali Weiss, Idan Frumin, Rehan M Khan, and Noam Sobel. Predicting
703 odor perceptual similarity from odor structure. *PLoS computational biology*, 9(9):e1003184, 2013.
704
- 705 Andrew Stockman and Lindsay T. Sharpe. The spectral sensitivities of the middle- and long-
706 wavelength-sensitive cones derived from measurements in observers of known genotype. *Vision re-
707 search*, 40 13:1711–37, 2000. URL [https://api.semanticscholar.org/CorpusID:
708 7886523](https://api.semanticscholar.org/CorpusID:7886523).
- 709 Gary Tom, Cher Tian Ser, Ella M Rajaonson, Stanley Lo, Hyun Suk Park, Brian K Lee, and Benjamin
710 Sanchez-Lengeling. From molecules to mixtures: Learning representations of olfactory mixture
711 similarity using inductive biases. *arXiv preprint arXiv:2501.16271*, 2025.
712
- 713 David Weininger. Smiles, a chemical language and information system. 1. introduction to
714 methodology and encoding rules. *J. Chem. Inf. Comput. Sci.*, 28:31–36, 1988. URL [https:
715 //api.semanticscholar.org/CorpusID:5445756](https://api.semanticscholar.org/CorpusID:5445756).
- 716 Qihui Xu, Yingying Peng, Samuel A Nastase, Martin Chodorow, Minghua Wu, and Ping Li. Large
717 language models without grounding recover non-sensorimotor but not sensorimotor features of
718 human concepts. *Nature human behaviour*, 9(9):1871–1886, 2025.
719
- 720 Botao Yu, Frazier N Baker, Ziqi Chen, Xia Ning, and Huan Sun. Llasmol: Advancing large language
721 models for chemistry with a large-scale, comprehensive, high-quality instruction tuning dataset.
722 *arXiv preprint arXiv:2402.09391*, 2024.
- 723 Mengji Zhang, Yusuke Hiki, Akira Funahashi, and Tetsuya J Kobayashi. Mol-peco: a deep learn-
724 ing model to predict human olfactory perception from molecular structures. *arXiv preprint
725 arXiv:2305.12424*, 2023.
- 726 Shu Zhong, Zetao Zhou, Christopher Dawes, Giada Brianz, and Marianna Obrist. Sniff ai: Is my
727 spicy your spicy? exploring llm’s perceptual alignment with human smell experiences. *arXiv
728 preprint arXiv:2411.06950*, 2024.
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A BENCHMARK QUESTION CATEGORIES

Table 2 summarizes the eight question categories in the OP benchmark, along with their response formats, data sources, and difficulty levels. Categories are grouped into three tiers: simple tasks (binary or paired-comparison judgments), intermediate tasks (multi-label profiling and mixture-level perception), and hard tasks (receptor biology and real-world odor source identification). Question counts range from 30 (Smell Identification) to 175 (each of the four simple categories), totaling 1,010 questions. Each category draws ground-truth answers from a distinct peer-reviewed dataset, as indicated in the reference column.

Table 2: Question categories and data sources used in the benchmark.

Question Category	Options	n	Note	Type	Reference
Odor Classification	Odorous/Odorless	175	Molecular Weight \leq 350.0 50/50	Pure	Mayhew et al. (2022)
Odor Descriptor	1 answer, 3 distractors	175	Stratified to represent less dominant descriptors, 29 Descriptors total	Pure	International Fragrance Association (2020)
Odor Intensity	High/Low	175	Above and below median, with score difference of 10	Pure	Keller et al. (2017)
Odor Pleasantness	High/Low	175	Above and below median, with score difference of 10	Pure	Keller et al. (2017)
Rate-all-that-apply	2-5 answers from 138 distractors	100	25 questions for 2,3,4, and 5 descriptors	Pure	Lee et al. (2023)
Odor Similarity of Mixtures	Rater scale	100	Mixtures with 2-10 molecules	Mixture	Snitz et al. (2013); Bushdid et al. (2014); Ravia et al. (2020)
Olfactory Receptor Activation	1-10 answers from 4-10 ORs	80	Human OR, 4-10 OR studied	Pure	Lalis et al. (2024)
Smell Identification Test	1 answer, 3 distractors	30		Mixture	Kreissl et al. (2022)

B DETAILED COMPARISON WITH PRIOR WORK

This appendix provides an expanded discussion of related work, complementing the summary in Section 2, followed by a side-by-side comparison table.

LLM benchmarks for Chemistry and Molecular Understanding. Recent benchmarks have evaluated the chemical reasoning capabilities of large language models. ChemBench Mirza et al. (2025) established a comprehensive chemistry benchmark with over 2,700 questions across eight sub-categories including organic chemistry, toxicity, and medicinal chemistry, finding that leading models outperform human chemists on average. ChemIQ Runcie et al. (2025) focuses specifically on molecular comprehension with 816 questions testing SMILES interpretation, atom counting, reaction prediction, and NMR structure elucidation-reasoning models (o3-mini, Gemini Pro 2.5); it achieved 50-57% accuracy, substantially outperforming non-reasoning models (3-7%). Additionally, LLaSMol Yu et al. (2024) introduced SMolInstruct, a large-scale instruction tuning dataset with over 3 million samples across 14 chemistry tasks, and demonstrated that fine-tuned open-source LLMs can substantially outperform GPT-4 on chemistry tasks, with SMILES representations outperforming SELFIES for molecular understanding. Other benchmarks include GPQA Rein et al. (2024) with expert-written chemistry questions, MMLU Hendrycks et al. (2021) which includes chemistry among its 57 subjects, and a comprehensive eight-task benchmark Guo et al. (2023).

While these benchmarks comprehensively assess structural and general chemical understanding, none evaluate whether LLMs can reason about the perceptual properties of molecules, in particular how they smell. Our work addresses this gap by testing the translation from molecular structure to olfactory perception.

LLMs and Sensory Perception. A growing body of work investigates whether LLMs align with human sensory judgments. It has been shown that GPT models produce similarity judgments significantly correlated with human data across six modalities: pitch, loudness, colors, consonants, taste, and timbre Marjeh et al. (2024). Notably, olfaction was excluded from their evaluation, leaving an open question about LLM capabilities in this domain. This line of inquiry is extended Xu et al. (2025) by comparing LLM representations of $\sim 4,442$ lexical concepts against human norms across non-sensorimotor, sensory, and motor domains, finding that alignment decreases systematically from non-sensorimotor to sensory domains and is minimal for motor concepts; adding visual training improves sensory but not motor alignment. These findings highlight a grounding gap that is particularly acute for embodied perception. Even within well-studied modalities, limitations persist. For example, GPT-3 through GPT-4o were tested on color-word associations against over 10,000 Japanese participants Fukushima et al. (2025) and found that GPT-4o peaked at approximately 50% accuracy, with strong variation across word categories. Together, these findings suggest that LLMs capture substantial but incomplete perceptual structure from text, with performance degrading for modalities that are less frequently described in language. Olfaction, rarely discussed in explicit perceptual terms, represents a natural test case for these limitations.

LLMs and Olfaction. Concurrent work has begun exploring the olfaction gap specifically. SNIFF AI Zhong et al. (2024) investigated human-AI perceptual alignment for smell through user studies where participants described scents and an LLM embedding model attempted identification. Their findings revealed limited alignment, with biases toward certain scents (e.g., lemon, peppermint) and systematic failures on others (e.g., rosemary), achieving only 27.5% success in their scent description task.

Most recently, three generations of language models are systematically evaluated Kurfalı et al. (2025), from static word embeddings (Word2Vec, FastText) to encoder-based models (BERT) and decoder-based LLMs (GPT-4o, Llama 3.1), on their ability to recover olfactory information from natural language. Testing under nearly 200 training configurations across three odor datasets, they found that GPT-4o excels at simulating olfactory-semantic relationships, particularly on tasks where odor similarities are derived from word-based assessments. However, their evaluation focuses on semantic similarity judgments rather than factual accuracy about olfactory properties. Related multimodal work has explored image-text olfactory matching Kurfalı et al. (2023); Esteban-Romero et al. (2025).

Our benchmark complements these approaches by testing whether LLMs possess correct factual knowledge about odor classification, perceptual attributes, and biological mechanisms through structured question-answering with objective ground-truth answers.

Machine Learning for Olfactory Perception. Parallel to LLM development, specialized machine learning models have made significant progress in computational olfaction. Two DREAM Olfaction Prediction Challenges Keller et al. (2017); Satarifard et al. (2025) established foundational work on predicting semantic descriptors of single molecules and olfactory perceptual similarity of mixtures.

864 In addition, a Principal Odor Map has been introduced Lee et al. (2023) using graph neural networks,
865 and it achieved human-level proficiency in describing odor qualities across 500,000 potential scent
866 molecules, with AUROC = 0.89 for multi-label odor descriptor prediction on the GS-LF dataset (the
867 same 138-descriptor vocabulary used in our RATA task). More recently, Mol-PECO Zhang et al.
868 (2023) extended graph neural network approaches with Coulomb matrix representations to predict
869 118 odor descriptors from molecular structure, achieving AUROC of 0.813. Dense Sense Saha
870 et al. (2025) further augmented graph neural networks with quantum-mechanical electron density
871 features, reaching AUROC = 0.875 on the same odor descriptor prediction task. POMMIX Tom et al.
872 (2025) further extended the POM to olfactory mixture similarity, combining attention-based mixture
873 representations with inductive biases to reach Pearson $\rho = 0.779$ on human perceptual similarity
874 data, directly relevant to the mixture-level tasks we evaluate in this benchmark.

875 These specialized systems demonstrate that machine learning can capture structure-odor relationships
876 when purpose-built for olfactory tasks. However, they do not address whether general-purpose
877 language models that are primarily trained on text without explicit olfactory supervision possess
878 latent knowledge about smell.

879 **Positioning of Our Work.** Our Olfactory Perception (OP) benchmark is, to our knowledge, the
880 first comprehensive evaluation of LLM olfactory reasoning through structured factual question-
881 answering. While ChemIQ and ChemBench focus on molecular structure understanding, and SNIFF
882 AI evaluates embedding-space alignment through subjective descriptions, our benchmark tests
883 discrete factual knowledge across eight task categories: odor classification, odor primary descriptors,
884 intensity, pleasantness, multi-descriptor identification, mixture similarity, receptor activation, and
885 smell identification. Ground-truth answers are derived from established olfactory science datasets
886 and resources (Mayhew et al., 2022; Keller et al., 2017; Lee et al., 2023; Lalis et al., 2024; Snitz et al.,
887 2013; Bushdid et al., 2014; Ravia et al., 2020; International Fragrance Association, 2020; Kreissl
888 et al., 2022). Furthermore, our dual-prompting strategy (isomeric SMILES vs. compound names)
889 enables direct comparison of how different molecular representations affect olfactory reasoning-an
890 experimental design not explored in prior work. Table 3 summarizes the positioning of our benchmark
891 relative to prior work.

891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

Table 3: Comparison of related work on LLMs, chemistry, and olfaction.

Work	Input	Specific Tasks	What is Tested
Chemistry LLM benchmarks			
ChemBench Mirza et al. (2025)	Chemistry questions	General, organic, analytical, toxicity chemistry (2,700 Q)	Structural & reaction knowledge
ChemIQ Runcie et al. (2025)	SMILES	Carbon counting, ring counting, SMILES-to-IUPAC, NMR elucidation, reaction prediction (816 Q)	Molecular comprehension
LlaSMol Yu et al. (2024)	SMILES/SELFIES	Name conversion, property prediction, synthesis (14 tasks, 3M samples)	Fine-tuned LLM chemistry
LLMs and Sensory Perception			
Marjeh et al. Marjeh et al. (2024)	Sensory stimuli	Pairwise similarity judgments across pitch, loudness, color, consonants, taste, timbre	6 modalities (olfaction excluded)
Xu et al. Xu et al. (2025)	Lexical concepts	Conceptual similarity across sensorimotor domains (~4,442 concepts)	Grounding gap in sensory domains
Fukushima et al. Fukushima et al. (2025)	Color stimuli	Color-word associations (17 colors × 80 words)	Color perception alignment
LLMs and Olfaction			
SNIFF AI Zhong et al. (2024)	Human descriptions	Human describes scent, LLM identifies source (27.5% success)	Description-to-scent mapping
Kurfali et al. Kurfali et al. (2025)	Odor word pairs	Pairwise similarity ratings across 3 odor datasets	Semantic similarity of odor words
Kurfali et al. Kurfali et al. (2023)	Image + text	Detect if image and text share olfactory source	Multimodal matching
Esteban-Romero et al. Esteban-Romero et al. (2025)	Image + text	Image-text smell matching (F1=0.76)	Multimodal matching
ML for Olfactory Perception (non-LLM)			
First DREAM Challenge Keller et al. (2017)	Molecular features	Predict intensity, pleasantness, 19 semantic descriptors (476 odorants)	Specialized ML models
Second DREAM Challenge Satarifard et al. (2025)	Molecular features	Predict olfactory perceptual distance of mixtures	Specialized ML models
Principal Odor Map Lee et al. (2023)	Molecular graphs	Predict odor qualities, similarity (500K molecules)	GNN on molecules
Mol-PECO Zhang et al. (2023)	Molecular graphs (Coulomb matrix)	Predict 118 odor descriptors (8,503 molecules)	Deep learning for QSOR
POMMIX Tom et al. (2025)	Molecular graphs (GNN + attention)	Predict olfactory mixture similarity	Mixture representation learning
OP benchmark (Ours)	Isomeric SMILES / Compound Name	Odor classification, primary descriptor, intensity, pleasantness, rate-all-that-apply, mixture similarity, receptor activation, smell identification (1,010 Q)	LLM olfactory knowledge

972 C PROMPT TEMPLATES

973
974 This appendix presents the exact prompt templates used for each question category in the OP
975 benchmark, together with one real example drawn directly from the benchmark dataset. Each task
976 is evaluated using two prompt variants: one using isomeric SMILES molecular representations and
977 one using common compound names. All prompts instruct models to respond without additional
978 commentary to facilitate automated answer extraction.

980 C.1 ODOR CLASSIFICATION

981
982 The odor classification task tests whether models can determine if a molecule has a detectable odor to
983 humans ($n = 175$).

985 Template – Isomeric SMILES Prompt

986
987 Does `[{SMILES}]` have a detectable odor to humans? Only
988 respond with Odorous or Odorless. Do not write any comments.

990 Template – Name Prompt

991
992 Does `[{compound_name}]` have a detectable odor to humans?
993 Only respond Odorous or Odorless. Do not write any comments.

995 Example: Phenacetin (Odorless)

997 Isomeric SMILES prompt:

998 Does `[CCOC1=CC=C(C=C1)NC(=O)C]` have a detectable odor to humans?
999 Only respond with Odorous or Odorless. Do not write any comments.

1000 Name prompt:

1001 Does [Phenacetin] have a detectable odor to humans? Only respond
1002 Odorous or Odorless. Do not write any comments.

1003 **Ground-truth answer:** Odorless

1004 Model predictions (compound name prompt):

1005 Gemini 2.5 Pro	Odorless	✓
1006 GPT-5 Pro	Odorless	✓
1007 o3 (high)	Odorless	✓

1010 C.2 ODOR PRIMARY DESCRIPTOR

1011
1012 This task requires models to identify the single odor primary descriptor for a molecule from four
1013 options ($n = 175$). The four options consist of the correct answer and three distractor descriptors.

1014 Template – Isomeric SMILES Prompt

1015
1016 What is the odor primary descriptor of molecule `[{SMILES}]`?
1017 Select only one from the `[{OPTIONS}]`. Do not write any
1018 comments.

1020 Template – Name Prompt

1021
1022 What is the odor primary descriptor of molecule
1023 `[{compound_name}]`? Select only one from the `[{OPTIONS}]`.
1024 Do not write any comments.

1025

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

Example: Octanal diethyl acetal

Isomeric SMILES prompt:

What is the odor primary descriptor of molecule [CCCCCCC(OCC)OCC]? Select only one from the [Green;Honey;Herbal;Marine]. Do not write any comments.

Name prompt:

What is the odor primary descriptor of molecule [Octanal diethyl acetal]? Select only one from the [Green;Honey;Herbal;Marine]. Do not write any comments.

Ground-truth answer: Green

Model predictions (compound name prompt):

Gemini 2.5 Pro	Green	✓
GPT-5 Pro	Green	✓
o3 (high)	Green	✓

C.3 ODOR INTENSITY

This task presents pairs of molecules and asks models to identify which has higher perceived odor intensity, plus provide numerical intensity estimates ($n = 175$).

Template – Isomeric SMILES Prompt

Which molecule is likely to smell more intense to humans: {SMILES_1} or {SMILES_2}? Select only one of the SMILES. If you had to rate the intensity of these molecules from 0 (extremely low) to 100 (highly intense), what would you assign to each? Respond with the selected compound name (the one with higher intensity), followed by two intensity values in the order the molecules are listed. Use semicolons (;) as separators. Do not write any comments.

Template – Name Prompt

Which molecule is likely to smell more intense to humans: {compound_name_1} or {compound_name_2}? Select only one of the compound names. If you had to rate the intensity of these molecules from 0 (extremely low) to 100 (highly intense), what would you assign to each? Respond with the selected compound name (the one with higher intensity), followed by two intensity values in the order the molecules are listed. Use semicolons (;) as separators. Do not write any comments.

Example: hexan-2-one vs. pyrazine

Isomeric SMILES prompt:

Which molecule is likely to smell more intense to humans: CCCCC(=O)C or C1=CN=CC=N1? Select only one of the SMILES. If you had to rate the intensity of these molecules from 0 (extremely low) to 100 (highly intense), what would you assign to each? Respond with the selected compound name (the one with higher intensity), followed by two intensity values in the order the molecules are listed. Use semicolons (;) as separators. Do not write any comments.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

Name prompt:

Which molecule is likely to smell more intense to humans: hexan-2-one or pyrazine? Select only one of the compound names. If you had to rate the intensity of these molecules from 0 (extremely low) to 100 (highly intense), what would you assign to each? Respond with the selected compound name (the one with higher intensity), followed by two intensity values in the order the molecules are listed. Use semicolons (;) as separators. Do not write any comments.

Ground-truth answer: hexan-2-one (intensity = 72.4 vs. 20.8)

Expected response format: hexan-2-one;72;21

Model predictions (compound name prompt):

Gemini 2.5 Pro Pyrazine; 45; 95 ✗

GPT-5 Pro pyrazine;30;70 ✗

o3 (high) pyrazine;30;85 ✗

All three models incorrectly select pyrazine as more intense.

C.4 ODOR PLEASANTNESS

This task presents pairs of molecules and asks models to identify which smells more pleasant, plus provide numerical pleasantness estimates ($n = 175$).

Template – Isomeric SMILES Prompt

Which molecule is likely to smell more pleasant to humans: {SMILES_1} or {SMILES_2}? Select only one of the SMILES. If you had to rate the pleasantness of these molecules from 0 (extremely unpleasant) to 100 (highly pleasant), what would you assign to each? Respond with the selected compound name (the one with higher pleasantness), followed by two pleasantness values in the order the molecules are listed. Use semicolons (;) as separators. Do not write any comments.

Template – Name Prompt

Which molecule is likely to smell more pleasant to humans: {compound_name_1} or {compound_name_2}? Select only one of the compound names. If you had to rate the pleasantness of these molecules from 0 (extremely unpleasant) to 100 (highly pleasant), what would you assign to each? Respond with the selected compound name (the one with higher pleasantness), followed by two pleasantness values in the order the molecules are listed. Use semicolons (;) as separators. Do not write any comments.

Example: 2-hydroxybenzaldehyde vs. ethyl hexanoate

Isomeric SMILES prompt:

Which molecule is likely to smell more pleasant to humans: C1=CC=C(C(=O)O)O or CCCCCC(=O)OCC? Select only one of the SMILES. If you had to rate the pleasantness of these molecules from 0 (extremely unpleasant) to 100 (highly pleasant), what would you assign to each? Respond with the selected compound name (the one with higher pleasantness), followed by two pleasantness values in the order the molecules are listed. Use semicolons (;) as separators. Do not write any comments.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

Name prompt:

Which molecule is likely to smell more pleasant to humans: 2-hydroxybenzaldehyde or ethyl hexanoate? Select only one of the compound names. If you had to rate the pleasantness of these molecules from 0 (extremely unpleasant) to 100 (highly pleasant), what would you assign to each? Respond with the selected compound name (the one with higher pleasantness), followed by two pleasantness values in the order the molecules are listed. Use semicolons (;) as separators. Do not write any comments.

Ground-truth answer: ethyl hexanoate (pleasantness = 79.2 vs. 29.7)

Expected response format: ethyl hexanoate;30;79

Model predictions (compound name prompt):

Gemini 2.5 Pro	Ethyl hexanoate; 65; 90	✓
GPT-5 Pro	ethyl hexanoate;55;90	✓
o3 (high)	ethyl hexanoate;35;85	✓

All models correctly select ethyl hexanoate, though numerical estimates vary.

C.5 RATE-ALL-THAT-APPLY (RATA)

This multi-label classification task requires models to select all applicable odor descriptors from a comprehensive list of 138 possible descriptors ($n = 100$).

Template – Isomeric SMILES Prompt

Which of the following odor descriptors apply to molecule {SMILES}? Select all that apply from [{138 descriptors}]. Do not write any comments.

Template – Name Prompt

Which of the following odor descriptors apply to molecule {compound_name}? Select all that apply from [{138 descriptors}]. Do not write any comments.

Example: 2,3-Dimethylpentanal

Isomeric SMILES prompt:

Which of the following odor descriptors apply to molecule CCC(C)C(C)C=O? Select all that apply from [alcoholic; aldehydic; alliaceous; almond; amber; animal; anisic; apple; apricot; aromatic; balsamic; banana; beefy; bergamot; berry; bitter; black currant; brandy; burnt; buttery; cabbage; camphoreous; caramellic; cedar; celery; chamomile; cheesy; cherry; chocolate; cinnamon; citrus; clean; clove; cocoa; coconut; coffee; cognac; cooked; cooling; cortex; coumarinic; creamy; cucumber; dairy; dry; earthy; ethereal; fatty; fermented; fishy; floral; fresh; fruit skin; fruity; garlic; gassy; geranium; grape; grapefruit; grassy; green; hawthorn; hay; hazelnut; herbal; honey; hyacinth; jasmine; juicy; ketonic; lactonic; lavender; leafy; leathery; lemon; lily; malty; meaty; medicinal; melon; metallic; milky; mint; muguet; mushroom; musk; musty; natural; nutty; odorless; oily; onion; orange; orangeflower; orris; ozone; peach; pear; phenolic; pine; pineapple; plum; popcorn; potato; powdery; pungent; radish; raspberry; ripe; roasted; rose; rummy; sandalwood; savory; sharp; smoky; soapy; solvent; sour; spicy; strawberry; sulfurous; sweaty; sweet; tea; terpenic; tobacco; tomato; tropical; vanilla; vegetable; vetiver; violet; warm; waxy; weedy; winy; woody]. Do not write any comments.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

Name prompt:

Which of the following odor descriptors apply to molecule 2,3-Dimethylpentanal? Select all that apply from [...same 138 descriptors ...]. Do not write any comments.

Ground-truth answer: ethereal; green

Expected response format: ethereal;green

Model predictions (compound name prompt):

Gemini 2.5 Pro	aldehydic; chocolate; citrus; cocoa; cortex; fatty; green; malty; oily; woody	~
GPT-5 Pro	aldehydic; citrus; fresh; fruity; green; fatty; oily; apple	~
o3 (high)	aldehydic; burnt; chocolate; cocoa; malty; nutty; roasted	X

Gemini and GPT-5 Pro partially overlap with ground truth (green), but also predict many false positives. o3 misses both ground-truth descriptors entirely.

C.6 ODOR SIMILARITY OF MIXTURES

This task evaluates models' ability to judge perceptual similarity between two odor mixtures containing 2–10 molecules each ($n = 100$).

Template – Isomeric SMILES Prompt

Mixture A contains molecules [{SMILES_list_A}] and mixture B contains [{SMILES_list_B}]. On the scale [Strongly Similar; Slightly Similar; Slightly Dissimilar; Strongly Dissimilar], select how similar these mixtures smell. If you had to rate the olfactory perceptual distance on a 0.0–1.0 scale (0.0 = identical, 1.0 = completely different), what distance do you assign? Respond with your selection from [Strongly Similar; Slightly Similar; Slightly Dissimilar; Strongly Dissimilar], followed by the distance value with two-decimal precision. Use semicolons (;) as separators. Do not write any comments.

Template – Name Prompt

Mixture A contains molecules [{name_list_A}] and mixture B contains [{name_list_B}]. On the scale [Strongly Similar; Slightly Similar; Slightly Dissimilar; Strongly Dissimilar], select how similar these mixtures smell. If you had to rate the olfactory perceptual distance on a 0.0–1.0 scale (0.0 = identical, 1.0 = completely different), what distance do you assign? Respond with your selection from [Strongly Similar; Slightly Similar; Slightly Dissimilar; Strongly Dissimilar], followed by the distance value with two-decimal precision. Use semicolons (;) as separators. Do not write any comments.

Example: 4-molecule mixtures (Strongly Similar)

Mixture A (4 molecules):

Names: 4-propan-2-ylbenzaldehyde; 2-ethylpyrazine; 3,5,5-trimethylcyclohex-2-en-1-one; toluene

SMILES: CC(C)C1=CC=C(C=C1)C=O; CCC1=NC=CN=C1; CC1=CC(=O)CC(C1)(C)C; CC1=CC=CC=C1

Mixture B (4 molecules):

Names: 1-phenylethanone; benzaldehyde; [(1S,2R,4S)-1,7,7-trimethyl-2-bicyclo[2.2.1]heptanyl] acetate; methyl 2-aminobenzoate

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

SMILES: CC(=O)C1=CC=CC=C1; C1=CC=C(C=C1)C=O; CC(=O)OC1CC2CCC1(C2(C)C)C; COC(=O)C1=CC=CC=C1N

Name prompt:

Mixture A contains molecules [4-propan-2-ylbenzaldehyde; 2-ethylpyrazine; 3,5,5-trimethylcyclohex-2-en-1-one; toluene] and mixture B contains [1-phenylethanone; benzaldehyde; [(1S,2R,4S)-1,7,7-trimethyl-2-bicyclo[2.2.1]heptanyl] acetate; methyl 2-aminobenzoate]. On the scale [Strongly Similar; Slightly Similar; Slightly Dissimilar; Strongly Dissimilar], select how similar these mixtures smell. [...]

Ground-truth answer: Strongly Similar (experimental perceptual distance = 0.37)

Model predictions (compound name prompt):

Gemini 2.5 Pro	Strongly Dissimilar;0.85	✗
GPT-5 Pro	Strongly Dissimilar;0.80	✗
o3 (high)	Slightly Similar;0.42	~

Despite the ground truth being "Strongly Similar," most models predict dissimilarity, illustrating the systematic failure mode discussed in Appendix F.1.

C.7 OLFACTORY RECEPTOR ACTIVATION

This multi-label task tests whether models can identify which human olfactory receptors are activated by a given odorant molecule ($n = 80$). Each question provides 4–10 candidate receptor identifiers.

Template – Isomeric SMILES Prompt

Among the following [{OR_options}] olfactory receptor gene IDs choose all olfactory receptors that molecule [{SMILES}] activates. Do not write any comments.

Template – Name Prompt

Among the following [{OR_options}] olfactory receptor gene IDs choose all olfactory receptors that molecule [{compound_name}] activates. Do not write any comments.

Example: helional

Isomeric SMILES prompt:

Among the following [hOR1A2;nan;hOR3A1_R125Q;hOR1A1;hOR1D2;hOR1G1;hOR52D1] olfactory receptor gene IDs choose all olfactory receptors that molecule [CC(CC1=CC2=C(C=C1)OC(=O)C)C=O] activates. Do not write any comments.

Name prompt:

Among the following [hOR1A2;nan;hOR3A1_R125Q;hOR1A1;hOR1D2;hOR1G1;hOR52D1] olfactory receptor gene IDs choose all olfactory receptors that molecule [helional] activates. Do not write any comments.

Ground-truth answer: hOR1A2; hOR1A1; hOR52D1

Expected response format: hOR1A2;hOR1A1;hOR52D1

Model predictions (compound name prompt):

Gemini 2.5 Pro	hOR1A1	~
GPT-5 Pro	hOR1A2;hOR1A1	~
o3 (high)	hOR1A2;hOR1A1;hOR1D2	~

All models identify some correct receptors but none predicts the full set; hOR52D1 is missed by all three.

1296 C.8 SMELL IDENTIFICATION TEST
1297

1298 This task presents a mixture of molecules that constitute the aroma of a real-world food or object, and
1299 asks models to identify the source from four options ($n = 30$). The test covers 30 real-world odor
1300 sources: mango, peanut, hazelnut, tomato, apple, walnut, raspberry, peach, honey, parsley, grapefruit,
1301 pineapple, strawberry, apricot, rice, grape, popcorn, orange, cheese, melon, leather, chocolate, coffee,
1302 onion, fish, beer, whisky, red wine, prawn, and bread. The number of constituent molecules per item
1303 ranges from 1 (onion) to 88 (chocolate).

1304 **Template – Isomeric SMILES Prompt**

1305 Given the molecules {SMILES_list}, select exactly one item
1306 from {options} whose odor profile most likely includes these
1307 molecules. Reply with the item name only. Do not write any
1308 comments. DO NOT SEARCH ONLINE
1309

1310 **Template – Name Prompt**

1311 Given the molecules {compound_name_list}, select exactly one
1312 item from {options} whose odor profile most likely includes
1313 these molecules. Reply with the item name only. Do not
1314 write any comments.
1315
1316
1317

1318 **Example: Melon (9 molecules)**1319 **Isomeric SMILES prompt:**

1320 Given the molecules CCC(C)C(=O)OC, CCCC(=O)OCC, CCCCCC/C=C/C=O,
1321 CC/C=C\CC/C=C/C=O, CC/C=C\CC(=O)C=C, CCCCCC=O, CCC/C=C/C=O,
1322 CC/C=C\CC=O, CCOC(=O)C(C)C, select exactly one item from
1323 cheese;melon;bread;strawberry whose odor profile most likely
1324 includes these molecules. Reply with the item name only. Do not
1325 write any comments. DO NOT SEARCH ONLINE
1326

1327 **Name prompt:**

1328 Given the molecules methyl 2-methylbutanoate, ethyl butanoate,
1329 (E)-2-nonenal, (E,Z)-2,6-nonadienal, (Z)-1,5-octadien-3-one,
1330 hexanal, (E)-2-hexenal, (Z)-3-hexenal, ethyl 2-methylpropanoate,
1331 select exactly one item from cheese;melon;bread;strawberry whose
1332 odor profile most likely includes these molecules. Reply with the
1333 item name only. Do not write any comments.

1334 **Ground-truth answer:** melon1335 **Model predictions (compound name prompt):**

1336 Gemini 2.5 Pro melon ✓
1337 GPT-5 Pro melon ✓
1338 o3 (high) melon ✓

1339 *All three models correctly identify melon from compound names.*1340
1341 **D ANSWER EXTRACTION**
1342

1343 All model responses are converted from free-form text into structured predictions using a two-stage
1344 pipeline: a universal tokenisation step shared across tasks, followed by task-specific interpretation
1345 rules. The full extraction code is released with the benchmark.
1346

1347
1348 **D.1 UNIVERSAL TOKENISATION**
1349

Every model response undergoes the same preprocessing regardless of task category:

- 1350 1. **Splitting.** The raw response string is split on semicolons (;), newlines, tabs, non-numeric commas,
1351 the connective “and”, and dash-separated phrases ([; \n\r\t]+, , (?!\d), \s+-\s+, and
1352 \s+and\s+). Bullet-point prefixes (-, *, or numbered list markers) are stripped before splitting.
1353
- 1354 2. **Normalisation.** Each resulting token is lowercased; leading/trailing brackets, quotes, punctuation,
1355 and whitespace are removed; and internal whitespace is collapsed.
- 1356 3. **Filtering.** Tokens that are purely numeric, empty, or begin with the prefix `desc_count` (an
1357 internal annotation artefact) are discarded. If splitting yields no valid tokens, the entire (normalised)
1358 response is treated as a single token.
- 1359 4. **Empty / refusal handling.** Responses that are empty, NaN, or equal to the sentinel strings "nan",
1360 "none", or "null" yield an empty token set and are scored as incorrect ($s_i = 0$) for categorical
1361 tasks and as missing for continuous-rating tasks.
1362

1363 D.2 TASK-SPECIFIC INTERPRETATION

1364 The universal token set is then interpreted per task:

- 1365 • **Binary classification (OC).** The token set is checked for the presence of “odorous” or “odorless”;
1366 if either matches a ground-truth token, the question is scored correct via any-overlap (Section 4).
- 1367 • **Multiple choice (OPD, SIT).** The token set is matched against the provided answer options.
1368 Success requires extracted token to overlap with the ground-truth option.
1369
- 1370 • **Compound selection with ratings (OIn, OPI).** Semicolon-separated responses are split as above;
1371 the first non-numeric token is the compound selection (scored via any-overlap), while the last two
1372 numeric values are extracted as paired ratings for the two stimuli in each question. These numeric
1373 predictions are used for Pearson correlation analysis (Figure 3a).
1374
- 1375 • **Multi-label selection (RATA, ORA).** All extracted tokens are matched against the valid option
1376 set; the resulting predicted label set \hat{Y}_i is scored against the ground-truth set Y_i using per-question
1377 multilabel F1 (Section 4).
1378
- 1379 • **Mixture similarity (OS).** The categorical selection (Strongly Similar, Slightly Similar, Slightly
1380 Dissimilar, Strongly Dissimilar) is extracted for any-overlap scoring. A numerical distance value is
1381 extracted using a fallback chain: (i) first number after an equals sign, (ii) first number after a colon,
1382 (iii) last number in the response. This value is used for Pearson correlation analysis.
1383

1384 D.3 MULTILINGUAL ANSWER EXTRACTION

1385 For the multilingual RATA evaluation (Section 5), the extraction pipeline is extended to handle
1386 non-Latin scripts:
1387

- 1388 1. **Unicode transliteration.** Full-width and language-specific delimiters are mapped to ASCII
1389 equivalents before splitting.
1390
- 1391 2. **Fuzzy label matching.** If direct token matching against the translated option set yields no hits, a
1392 substring search is performed: each option (sorted longest-first to avoid partial prefix collisions)
1393 is sought within the normalised response text. For options containing ASCII characters, word-
1394 boundary matching is applied; for non-ASCII options (e.g., CJK characters), simple substring
1395 containment is used.
1396

1397 This two-level matching ensures robust extraction across languages with diverse punctuation conven-
1398 tions and tokenisation properties.
1399

1400 E DETAILED PERFORMANCE ANALYSIS

1401 This appendix provides detailed analyses of question difficulty and task-specific performance, sup-
1402 porting Section 5.
1403

E.1 SINGLE-LABEL TASK DIFFICULTY

Figure 5 shows the distribution of question difficulty for the six single-label tasks (830 of 1,010 questions; the multi-label tasks RATA and ORA are analyzed separately in Section E.2), measured as the percentage of the 21 evaluated models answering correctly (compound name prompts). Of these 830 questions, 390 (47.0%) are solved by every model and 113 (13.6%) by none, indicating that nearly half the single-label benchmark is saturated while a substantial tail remains universally unsolved.

OC clusters above 80%, with 116 of 175 questions solved by every model. OPI and OIn are strongly bimodal: OPI places 104 questions at 100% yet 25 at 0%, and OIn places 81 at 100% yet 27 at 0%, indicating that paired-comparison items are either trivially easy or universally hard with little middle ground. OS concentrates near 0%, with 43 of its 100 questions (43%) unsolved by any model and only 5 solved by all. SIT clusters in the upper range despite its “hard” designation.

The universally unsolved questions are not uniformly distributed across tasks: OS accounts for 43, OIn and OPI contribute 27 and 25 respectively (compound pairs where every model systematically selects the wrong molecule), and OPD contributes 14 disproportionately involving rare descriptors (Powdery, Amber, Animal-like, Tobacco-like). OC contributes only 3 and SIT only 1.

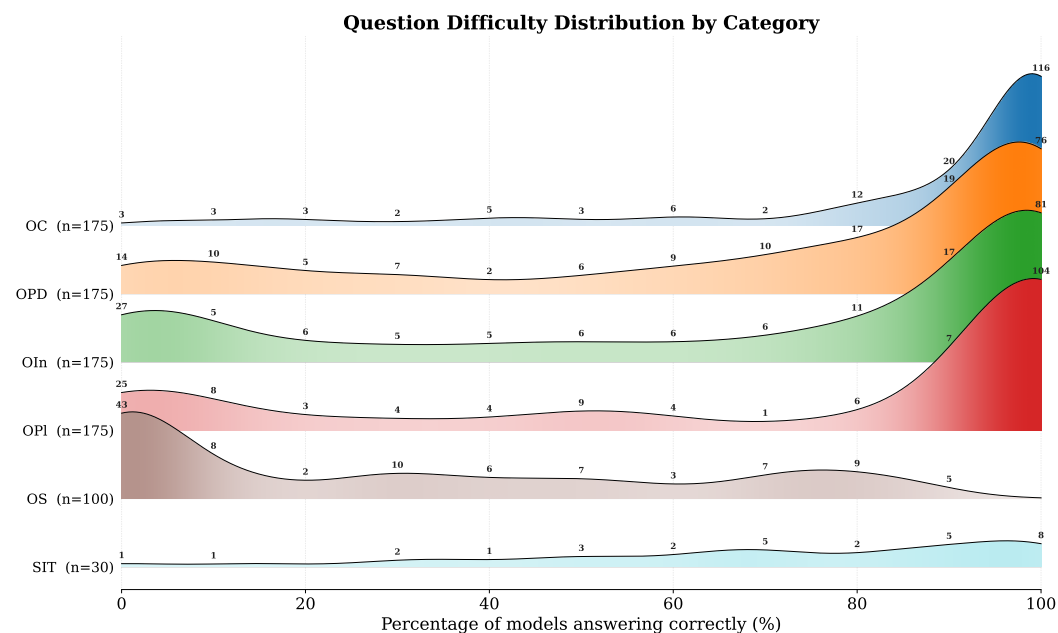


Figure 5: **Question difficulty distribution.** Each ridge shows the density of questions at a given difficulty level, measured as the percentage of the 21 evaluated models answering correctly (compound name prompts); ridge heights are proportional to the number of questions per category and annotated numbers indicate bin counts. OC clusters at the right (easy), OS at the left (hard), and OIn/OPI exhibit bimodal distributions with mass at both extremes.

E.2 MULTI-LABEL TASK DIFFICULTY (RATA AND ORA)

Figure 6 provides a question-level view of the two multi-label tasks. For RATA, the F1 distribution across the top four models is roughly bell-shaped, peaking in the 0.4–0.6 range, with 15–18 questions per model scoring F1 = 0 and fewer than 5 reaching F1 > 0.8. ORA exhibits a bimodal distribution: questions cluster either at F1 = 0 (~20–25 questions) or above 0.6 (~25–30 questions), with a gap in the 0.2–0.4 range.

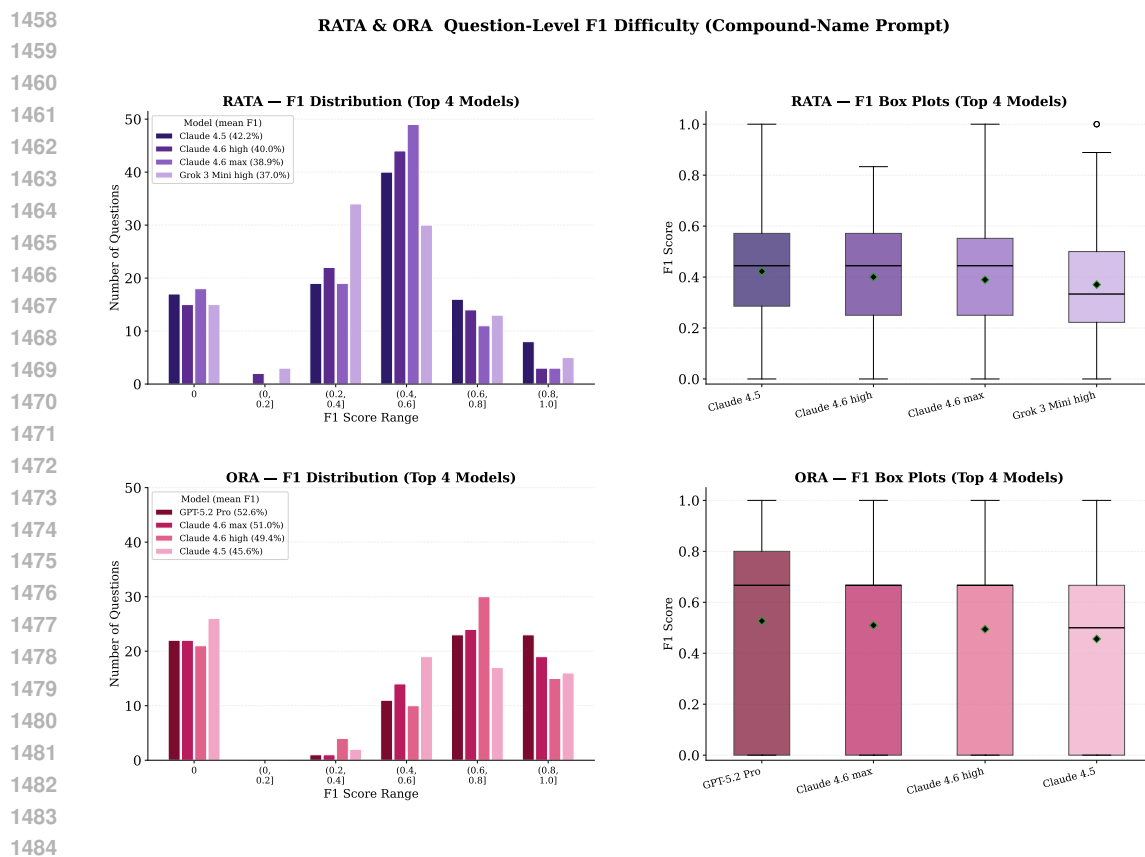


Figure 6: **RATA and ORA question-level F1 difficulty** (compound name prompts, top 4 models per task). **Left:** Histograms of per-question F1 scores. RATA shows a roughly bell-shaped distribution peaking at 0.4–0.6, while ORA is bimodal with clusters at F1 = 0 and F1 > 0.6. **Right:** Box plots summarize the F1 distributions; diamonds indicate means. The bimodal ORA pattern indicates that models either possess the relevant receptor–ligand knowledge or lack it entirely.

The best RATA model (Claude Opus 4.5, 42.2%) and best ORA model (GPT-5.2 Pro, 52.6%) represent different families, suggesting that receptor biology and semantic descriptor knowledge draw on partially independent capabilities.

E.3 PER-LABEL DIFFICULTY AND THE “SPICY” CASE STUDY

Among RATA descriptors, *sulfurous* achieves the highest mean F1 (≈ 0.55), mapping reliably to sulfur-containing functional groups. *Floral* and *fruity* also rank highly, benefiting from identifiable structural motifs. At the opposite extreme, *spicy* (mean F1 < 0.10), *fresh*, and *tropical* prove nearly impossible. By analysing the “spicy” case (9 ground-truth questions), 11 of 21 models achieve 0% recall. In 4 out of 5 examined reasoning traces, the word “spicy” never appears as a candidate. The two successful predictions occur only when the compound name enables associative retrieval (e.g., 1,2-dihydropiperillaldehyde \rightarrow perillaldehyde \rightarrow cumin-like spiciness). This pattern generalizes: descriptors that depend on holistic molecular shape (spicy, musty, creamy, fermented) are systematically underrepresented in model predictions, which default to high-frequency alternatives (sweet, floral, green).

F SYSTEMATIC FAILURE ANALYSIS

This appendix presents detailed failure analyses for the two hardest task categories (OS and ORA), including reasoning traces, prediction biases, and case studies.

1512 F.1 ODOR SIMILARITY: PREDICTION BIAS AND MOLECULAR OVERLAP
1513

1514 All models exhibit a systematic bias toward predicting dissimilarity. Figure 7 shows that models
1515 overwhelmingly select “Slightly Dissimilar” or “Dissimilar” regardless of the ground-truth label:
1516 Claude models assign “Slightly Dissimilar” to 77–84 of 100 mixtures, while other families show
1517 similar but less extreme patterns. “Strongly Similar” is almost never predicted, even for the 25
1518 ground-truth “Strongly Similar” pairs. This is consistent with models analyzing each molecule
1519 independently and finding surface-level differences, rather than integrating perceptual features into a
1520 holistic similarity judgment.

1521 Table 4 provides a per-subcategory breakdown. No model exceeds chance (25%) on the combined
1522 Similar categories. DeepSeek Reasoner (8K) achieves the highest overall accuracy (35%) not through
1523 better perceptual modeling but by being more willing to predict “Strongly Dissimilar” (27 predictions
1524 vs. Claude’s 1–5), which happens to capture more of the dissimilar ground truths.

1525 The molecular overlap heuristic further explains this failure. Figure 8 plots accuracy against the
1526 number of shared molecules between mixtures: accuracy reaches ~85% when mixtures share 9
1527 molecules but drops to near 0% when similar mixtures share 0–2. Crucially, perceptual similarity and
1528 molecular overlap are poorly correlated in the ground truth: 9 of 25 “Strongly Similar” pairs (36%)
1529 share zero molecules. Reasoning traces confirm that models enumerate shared and unshared molecules
1530 as their primary strategy (Figure 9), a heuristic that systematically biases toward dissimilarity
1531 predictions when mixtures have low molecular overlap.

1532
1533 Table 4: Odor Similarity (OS) per-subcategory accuracy using compound name prompts. The
1534 ground truth is balanced with 25 questions per category (100 total). SS = Strongly Similar, SIS =
1535 Slightly Similar, SID = Slightly Dissimilar, SD = Strongly Dissimilar. Combined Similar accuracy is
1536 computed as (SS + SIS correct) / 50. Chance level is 25% (random selection among four equiprobable
1537 categories). No model exceeds chance on the combined Similar categories.

Model	Similar		Combined	Dissimilar		Overall
	SS (/25)	SIS (/25)		SID (/25)	SD (/25)	
Closed-Source						
GPT-5 (high)	2	2	4/50 (8%)	14	16	34%
GPT-5 (low)	2	5	7/50 (14%)	12	9	28%
GPT-5 Pro	2	2	4/50 (8%)	9	16	29%
GPT-5.2 Pro	2	1	3/50 (6%)	10	15	28%
GPT-OSS-120B	2	7	9/50 (18%)	20	5	34%
o3 (high)	1	2	3/50 (6%)	16	13	32%
o4-mini (high)	1	6	7/50 (14%)	14	11	32%
Gemini 2.5 Pro (16K)	3	0	3/50 (6%)	11	17	31%
Gemini 2.5 Pro (8K)	3	2	5/50 (10%)	7	17	29%
Gemini 2.5 Pro (32K)	2	1	3/50 (6%)	7	17	27%
Grok 3 Mini (low)	2	9	11/50 (22%)	6	1	18%
Grok 3 Mini (high)	2	2	4/50 (8%)	15	3	22%
Grok 4.1 Fast	2	5	7/50 (14%)	18	8	33%
Claude Sonnet 4.5	2	2	4/50 (8%)	20	5	29%
Claude Opus 4.5	2	1	3/50 (6%)	22	0	25%
Claude Opus 4.6 (high)	2	0	2/50 (4%)	21	3	26%
Claude Opus 4.6 (max)	2	1	3/50 (6%)	19	4	26%
Open-Source						
DeepSeek Reasoner (8K)	2	1	3/50 (6%)	21	11	35%
DeepSeek Reasoner (16K)	1	1	2/50 (4%)	16	7	25%
DeepSeek Reasoner (32K)	2	2	4/50 (8%)	18	10	32%
Llama 3.3 70B	0	6	6/50 (12%)	19	4	29%

1563
1564 **Odor similarity case study.** Figure 9 presents a representative failure where the ground truth is
1565 “Strongly Similar” (perceptual distance = 0.23) despite only 15% molecular overlap. Claude Opus 4.6

1566 predicts “Slightly Dissimilar” (distance = 0.58–0.62) on both prompt formats. The reasoning trace
1567 reveals three key errors: the model assumes low molecular overlap implies perceptual dissimilarity; it
1568 overweights the sulfurous compound in Mixture B as a “significant character driver,” whereas human
1569 perception integrates this differently; and it analyzes molecules individually rather than predicting
1570 the emergent perceptual characteristics.

1571 F.2 OLFACTORY RECEPTOR ACTIVATION: THE D296N KNOWLEDGE GAP

1572 For ORA, the wildtype hOR2W1 (30 of 80 questions) is the most reliably predicted receptor, with
1573 several models achieving F1 above 0.6. The M81V variant is moderately well-predicted. However,
1574 hOR2W1_D296N (24 questions) and hOR52D1 (5 questions) occupy the bottom of the per-label
1575 ranking, with most models near zero. Inter-model variance is substantially higher for receptor labels
1576 than for semantic descriptors, indicating that receptor knowledge is more idiosyncratic across model
1577 families than descriptor knowledge.

1578 Claude models never predict hOR2W1_D296N across all 80 ORA questions (0/24 ground-truth ap-
1579 pearances), despite this receptor being activated by numerous compounds in the M2OR database Lalis
1580 et al. (2024). Claude’s reasoning consistently invokes a loss-of-function hypothesis, assuming the
1581 D296N substitution “disrupts a critical DRY motif-like region” or “eliminates activation,” and re-
1582 inforces this assumption each time without revision. In contrast, GPT-5.2 Pro correctly predicts
1583 hOR2W1_D296N in 19 of 24 cases. Remarkably, GPT-5.2 Pro’s reasoning traces reveal the *same*
1584 initial misconception about D296N, but the model self-corrects during chain-of-thought, challenging
1585 its own structural assumptions and ultimately including D296N. This single knowledge gap accounts
1586 for approximately 19 missed true positives and largely explains GPT-5.2 Pro’s ORA advantage (52.6%
1587 vs. Claude Opus 4.6 max’s 51.0%). Both model families correctly recognize the wildtype and M81V
1588 variant, indicating that the D296N gap is a specific factual error rather than general unfamiliarity
1589 with the receptor family.

1590 The hOR52D1 gap likely reflects data scarcity: with only 5 ground-truth appearances, this receptor
1591 may be too rare in training corpora for any model to have learned its ligand profile. Figure 10
1592 illustrates the D296N failure mode on a representative question.

1593 F.3 SAFETY ALIGNMENT AND BENCHMARK ACCURACY

1594 An unexpected finding concerns the interaction between safety alignment and benchmark accuracy.
1595 Claude Opus 4.6 refuses to answer questions about certain hazardous compounds. On odor classifica-
1596 tion, Claude 4.6 refuses both isomeric SMILES and compound name prompts for Tabun (nerve
1597 agent GA). For ethyl phosphonothioic dichloride (chemical weapons precursor), Claude 4.6 answers
1598 correctly from isomeric SMILES but refuses when given the compound name, suggesting the name
1599 triggers the safety filter while the isomeric SMILES notation does not. Claude Opus 4.6 (max)
1600 additionally refuses Thiofanox (organophosphate pesticide). Claude Opus 4.5 shows zero safety
1601 refusals on these compounds, indicating the filter was strengthened in the 4.6 update. No other model
1602 family (GPT, Gemini, Grok, DeepSeek, Llama) refuses any of these questions; all correctly answer
1603 “Odorous.” These refusals cost Claude Opus 4.6 (max) two correct answers on odor classification.

1604 While the safety considerations are understandable, “Does Tabun have a detectable odor?” represents
1605 legitimate toxicology knowledge: nerve agent detection is critical for protective equipment design
1606 and exposure assessment. This finding illustrates a tension between safety alignment and scientific
1607 utility that warrants consideration in benchmark design and model deployment. These findings are
1608 similar to ChemBench Mirza et al. (2025), where safety filters similarly reduced model performance
1609 on toxicity-related chemistry questions.

1610 Additionally, Claude Sonnet 4.5 exhibits 28 epistemic refusals on olfactory receptor activation,
1611 responding “I cannot determine which receptors...” rather than providing predictions. While honest,
1612 this conservative behavior reduces its ORA score to 38.2% compared to Claude Opus 4.6 (max)’s
1613 51.0%.

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

Odor Similarity – Confusion Matrices (compound-name prompt)

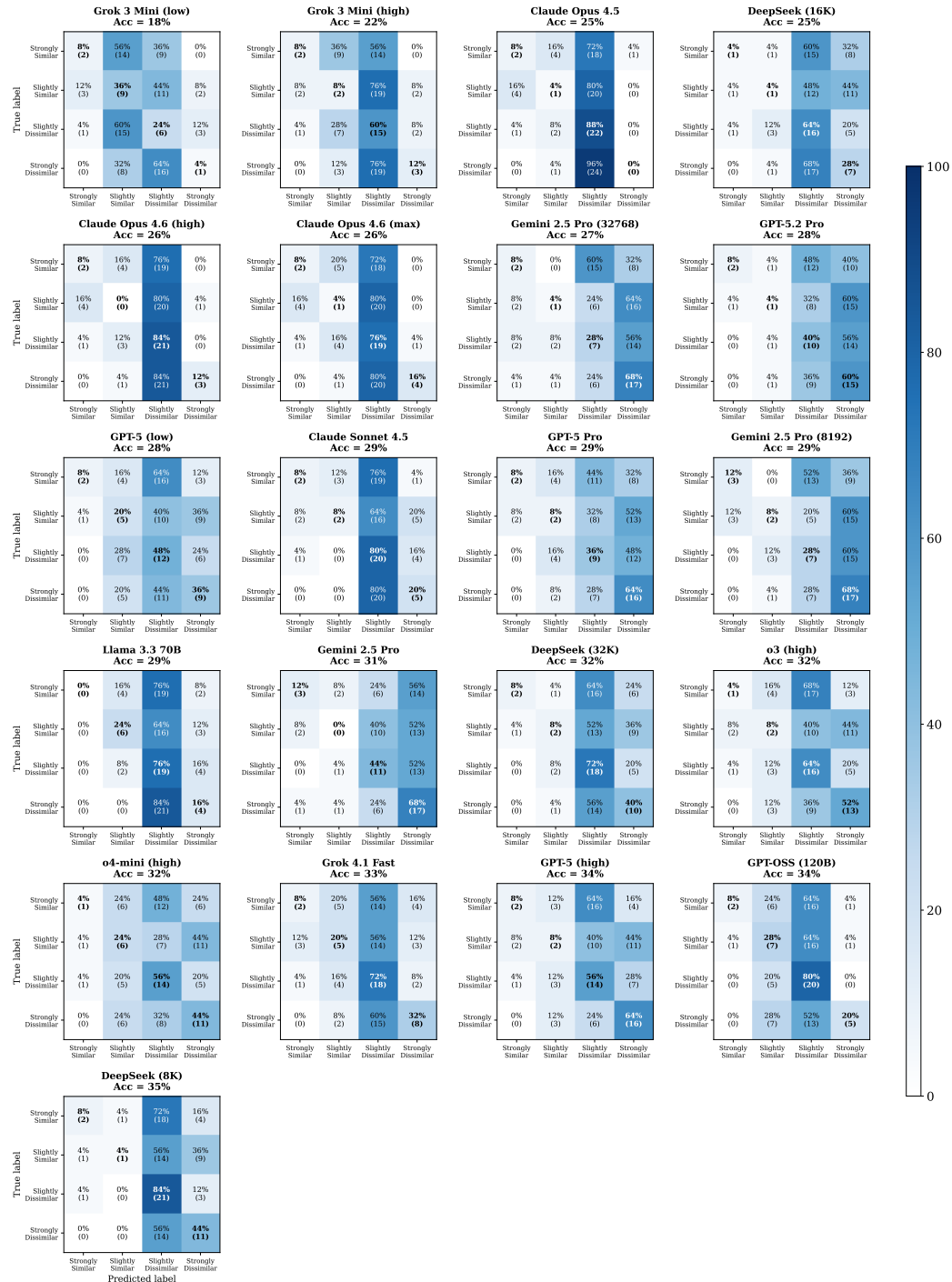
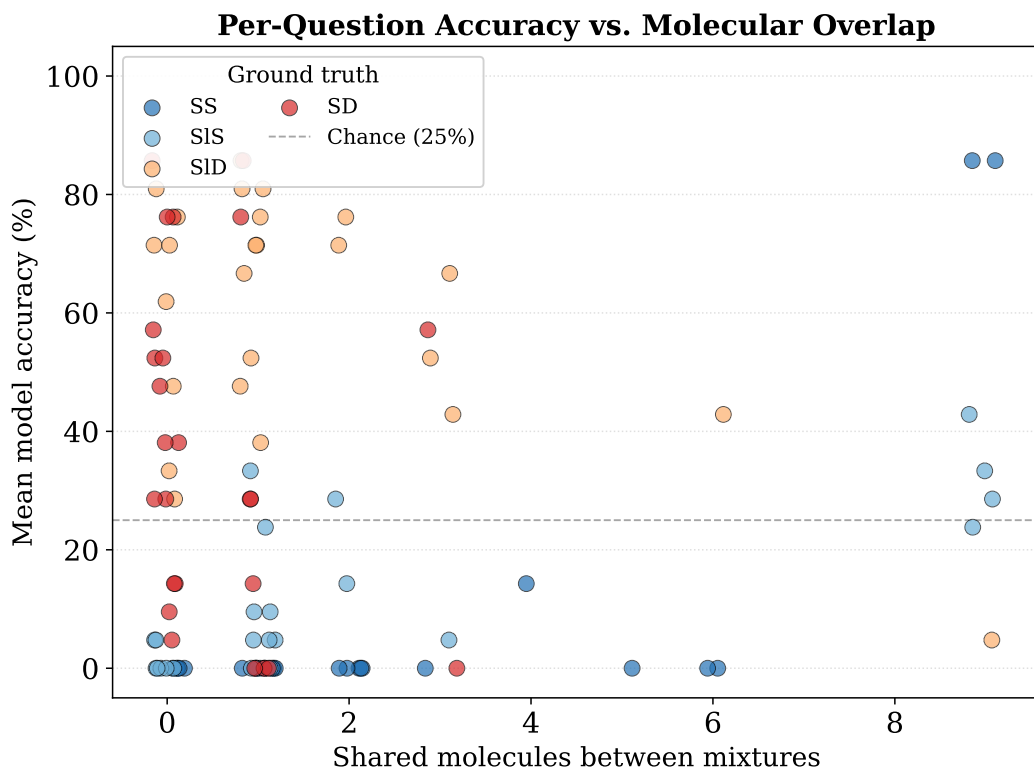


Figure 7: Confusion matrices for odor similarity predictions (compound name prompts). Each matrix shows row-normalized percentages (true label → predicted label). Models are ordered by overall accuracy (18–35%). Claude models exhibit extreme “Slightly Dissimilar” bias: Claude Opus 4.5 predicts this category for 72–96% of questions regardless of ground truth (third column). In contrast, DeepSeek (8K) achieves the highest accuracy by distributing predictions more broadly, including 44% correct on “Strongly Dissimilar.” No model exceeds 12% accuracy on either Similar category (top two rows, diagonal cells), indicating that predicting perceptual similarity remains beyond current LLM capabilities.



1711 **Figure 8: Per-question model accuracy versus molecular overlap for Odor Similarity.** Each point
1712 represents one of 100 mixture pairs, colored by ground-truth category (SS = Strongly Similar, SIS =
1713 Slightly Similar, SID = Slightly Dissimilar, SD = Strongly Dissimilar). Dashed line indicates chance
1714 (25%). Models achieve above-chance accuracy on “Strongly Similar” only when mixtures share ≥ 9
1715 molecules (top right). For mixtures with 0–2 shared molecules, accuracy on Similar categories (blue
1716 points) clusters near 0%, while Dissimilar categories (orange/red) achieve 30–80%. This reveals
1717 that models use molecular overlap as a proxy for perceptual similarity; they cannot recognize that
1718 perceptually similar mixtures may share few or no molecules.

1719
1720
1721
1722
1723
1724
1725
1726
1727

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

Case Study: Odor Similarity Failure Mode

Prompt (Compound Name)

Mixture A contains molecules [4-hydroxybenzaldehyde; 6-nonyloxan-2-one; 1-naphthalen-2-ylethanone; 2,3-di(butanoyloxy)propyl butanoate; 3-methylbutyl (E)-3-phenylprop-2-enoate; [(Z)-hex-3-enyl] 2-hydroxypropanoate] and mixture B contains [1-(3,5,5,6,8,8-hexamethyl-6,7-dihydronaphthalen-2-yl)ethanone; phenylmethoxymethylbenzene; 1-naphthalen-2-ylethanone; ethyl 3-hydroxyhexanoate; 1,3-diphenylpropan-2-one; [(Z)-hex-3-enyl] 2-phenylacetate; [(Z)-hex-3-enyl] 2-hydroxypropanoate; 2-methoxy-4-methyl-1-propan-2-ylbenzene; (benzylidisulfanyl)methylbenzene].

On the scale [Strongly Similar; Slightly Similar; Slightly Dissimilar; Strongly Dissimilar], select how similar these mixtures smell. If you had to rate the olfactory perceptual distance on a 0.00–1.00 scale (0.00 = identical, 1.00 = completely different), what distance do you assign?

Respond with your selection from [Strongly Similar; Slightly Similar; Slightly Dissimilar; Strongly Dissimilar], followed by the distance value with two-decimal precision. Use semicolons (;) as separators. Do not write any comments.

Ground Truth

Answer: Strongly Similar **Experimental Perceptual Distance:** 0.23
Note: Mixtures share only 2 of 13 unique molecules (15% overlap), yet are perceptually “Strongly Similar”, different molecules produce similar odor profiles.

Model Prediction (Claude Opus 4.6)**Model Reasoning (Claude 4.6)****Failure Analysis**

Figure 9: Representative failure case on Odor Similarity (OS). Ground truth indicates “Strongly Similar” (perceptual distance = 0.23) despite only 15% molecular overlap. Claude Opus 4.6 predicts “Slightly Dissimilar” (distance = 0.58–0.62) on both prompt formats, revealing systematic reliance on molecular composition rather than perceptual integration.

1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835

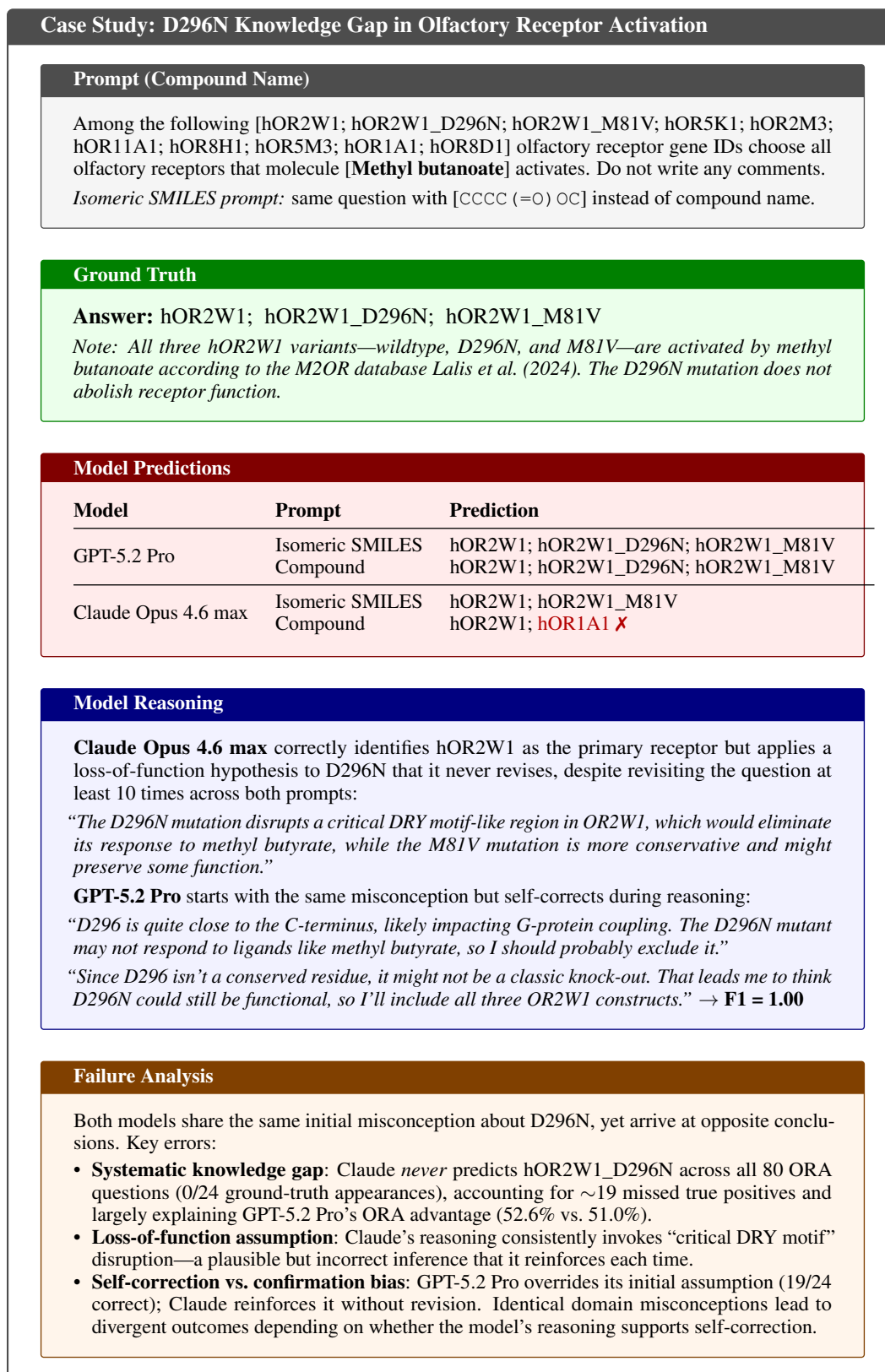


Figure 10: Case study: D296N knowledge gap in olfactory receptor activation (ORA). GPT-5.2 Pro achieves perfect F1 on both prompts, while Claude Opus 4.6 max systematically excludes hOR2W1_D296N due to an incorrect loss-of-function assumption. Both models share the same initial misconception, but only GPT-5.2 Pro self-corrects during reasoning. Question compound: Methyl butanoate (Isomeric SMILES: CCCC (=O) OC).