

Supplementary Information for: Thinking, Fast and Slow in AI Agents for Personalized Decision-making

Yang Zhao^{1, 2}, Yichen Lin³, Takeru Igusa^{1, 2}, and Hao Frank Yang^{1, 2, 4 *}

¹Department of Civil and Systems Engineering, Johns Hopkins University, Baltimore, MD, USA

²Center for Systems Science and Engineering, Johns Hopkins University, Baltimore, MD, USA

³Department of Applied Mathematics and Statistics, Johns Hopkins University, Baltimore, MD, USA

⁴Johns Hopkins Data Science and AI Institute, Johns Hopkins University, Baltimore, MD, USA

*Corresponding Author Email: haofrankyang@jhu.edu

1	Experiment Settings	2
1.1	Data Splits	2
1.2	Group Splits	2
2	Initial Utility Functions and Knowledge Sets	3
2.1	Initial Knowledge Sets	3
2.2	Initial Utility Functions	4
3	Algorithm Details	5
4	Example of Co-Gradients	6
5	Additional Results	7
5.1	Contributions of Grouped Features Settings • Results	7
5.2	Simplicity of Refined Utility Functions	9
5.3	Feature Evolution in Utility Functions	9
5.4	Cost Analysis Inference • Training	14
5.5	Optimized Utility Functions and Knowledge Sets Examples	15
5.6	Utility-Knowledge Contributions for All Groups and Tasks VAC • TMC • JOB • MAJ • PR • ALC	27
6	Ablation Study	45

1 Experiment Settings

1.1 Data Splits

Supplementary Table 1: Case counts in the training set, validation + testing set, and the total sample across six datasets. Data availability varies across demographic groups, with some groups containing substantially more observations than others. Consequently, the absolute case counts and their relative proportions may vary within a certain range across datasets and splits.

	Vaccine	Swissmetro	Job	Marijuana	Pain Reliever	Alcohol
Training	595	450	405	450	367	401
Validation + Testing	405	558	471	530	536	566
Total	1,000	1,008	876	980	903	967

1.2 Group Splits

Supplementary Table 2: Demographic group definitions across datasets. Income is reported in thousands of local currency units.

Group	Vaccine	Travel Mode	Job	Marijuana	Pain Reliever	Alcohol
Age						
Older	≥ 60	≥ 65	≥ 65	≥ 65	≥ 65	≥ 65
Middle-aged	39–60	40–65	35–65	35–65	35–65	35–65
Young	< 38	< 40	< 35	< 35	< 35	< 35
Income (in thousands)						
High	\geq county median	≥ 100 (CHF)	≥ 100 (\$)	≥ 75 (\$)	≥ 75 (\$)	≥ 75 (\$)
Middle	N/A	50–100 (CHF)	40–100 (\$)	50–75 (\$)	50–75 (\$)	50–75 (\$)
Low	$<$ county median	< 50 (CHF)	< 40 (\$)	< 50 (\$)	< 50 (\$)	< 50 (\$)

2 Initial Utility Functions and Knowledge Sets

2.1 Initial Knowledge Sets

Prompt 2.1: The initial knowledge set for the VAC task

- Belief that available COVID-19 vaccines are safe is the single biggest trigger for starting vaccination.
- Trust in science and government plays a critical role in vaccine uptake.

Prompt 2.2: The initial knowledge set for the TMC task

- Interaction between ticket price and income level significantly drives mode choice; higher-income travelers value time more, while lower-income groups are more price sensitive.
- Number of luggage, car availability, and annual pass ownership behind Swissmetro plays a critical role in uptake.
- If `has_ga_travel_pass` is 1, the ticket price is the price for travel pass. Otherwise, the ticket price is the price for single trip for swissmetro and train.

Prompt 2.3: The initial knowledge set for the JOB task

- Personal and family covid diagnose state is a factor when people make decision about whether to work
- Health or mental issues can effect work decision
- If have work from home requirement, people are more likely to work for salary or not work

Prompt 2.4: The initial knowledge set for the MAJ task

- Earlier `age_first_use` can establish more stable behavioral patterns, which may shorten the interval since the last occurrence.
- Lower perceived `difficulty_getting_marijuana` generally corresponds to more recent behavioral engagement because easier access reduces gaps between occurrences.
- Strong `urge_marijuana` reflects stronger impulses, which can be associated with shorter recency intervals compared to individuals without such impulses.

Prompt 2.5: The initial knowledge set for the ALC task

- Ever `binge_drinking` may reflect a more established behavioral pattern, which can shorten the interval since the last alcohol misuse.
- Lower perceived `risk_harm_5plus_drinks_weekly` can correspond to shorter intervals since the last alcohol misuse, because lower perceived harm tends to support more frequent engagement.
- Frequency `doing_risky_things` captures stronger behavioral impulses, which can be associated with shorter recency intervals compared to individuals with lower risk-taking tendencies.

Prompt 2.6: The initial knowledge set for the PR task

- Painreliever_disorder_past_year suggests stronger or more persistent engagement with the behavior, which can correspond to shorter intervals since the last use.
- Someone_selling_drugs_last_month indicates greater environmental exposure, potentially making more recent occurrences more likely.
- Times_attend_religious_services reflects structured routines and constraints, which can correspond to longer intervals since the last use.

2.2 Initial Utility Functions

VAC.

$$\begin{aligned}U_{\text{Unvaccinated}} &= C_1 \cdot \text{covid_threat} + C_2 \cdot \text{covid_preventable_by_vax} + C_3 \cdot \log(1 + \text{trust_science}) \\U_{\text{Booster}} &= C_1 \cdot \text{covid_threat} + C_2 \cdot \text{covid_preventable_by_vax} + C_3 \cdot (\text{nurse} + \text{healthcare_worker} + \text{physician}) \\U_{\text{Vaccinate}} &= C_1 \cdot \text{covid_threat} + C_2 \cdot \text{covid_preventable_by_vax} + C_3 \cdot \text{risk_of_covid_greater_than_vax}\end{aligned}\quad (1)$$

TMC.

$$\begin{aligned}U_{\text{Car}} &= C_1 \cdot \text{trip_purpose} + C_2 \cdot \text{annual_income_level} + C_3 \cdot \text{is_car_available} \\U_{\text{Train}} &= C_1 \cdot \text{trip_purpose} + C_2 \cdot \text{annual_income_level} + C_3 \cdot \text{ticket_payer_type} \\U_{\text{Swissmetro}} &= C_1 \cdot \text{trip_purpose} + C_2 \cdot \text{annual_income_level} + C_3 \cdot \text{has_ga_travel_pass}\end{aligned}\quad (2)$$

JOB.

$$\begin{aligned}U_{\text{Work_for_salary}} &= C_1 \cdot \log(\text{income} + 1) + C_2 \cdot \text{covid_worked_from_home} \\U_{\text{Self_employed}} &= C_1 \cdot \text{covid_worked_from_home} + C_2 \cdot \log(\text{age} + 1) \\U_{\text{Did_not_work}} &= C_1 \cdot \text{self_rated_health} + C_2 \cdot \text{tried_received_gov_assistance} \\&\quad + C_3 \cdot \text{restriction_work_from_home_req}\end{aligned}\quad (3)$$

MAJ.

$$\begin{aligned}U_{\text{Past_30_Days}} &= C_1 \cdot \text{strong_urge_marijuana} + C_2 \cdot \exp(-\text{difficulty_getting_marijuana}) \\&\quad + C_3 \cdot \log(1 + \text{age_first_use}) \\U_{\text{Past_12_Months}} &= C_1 \cdot \text{tried_to_stop_marijuana} + C_2 \cdot \text{difficulty_getting_marijuana} + C_3 \cdot \log(1 + \text{age}) \\U_{\text{More_Than_12_Months}} &= C_1 \cdot \text{risk_mj_weekly} + C_2 \cdot \text{opinion_adult_try_marijuana} + C_3 \cdot \exp(-\text{education_level})\end{aligned}\quad (4)$$

ALC.

$$\begin{aligned}U_{\text{Past_30_Days}} &= C_1 \cdot \text{ever_binge_drinking} + C_2 \cdot \text{frequency_doing_risky_things} \\U_{\text{Past_12_Months}} &= C_1 \cdot \text{risk_harm_5plus_drinks_weekly} + C_2 \cdot \text{workplace_drug_alcohol_policy} \\&\quad + C_3 \cdot \log(1 + \text{doctor_visits_last_year}) \\U_{\text{More_Than_12_Months}} &= C_1 \cdot \text{risk_harm_5plus_drinks_weekly} + C_2 \cdot \text{self_evalute_mental_health} \\&\quad + C_3 \cdot \exp(-\text{education_level})\end{aligned}\quad (5)$$

PR.

$$\begin{aligned}U_{\text{Past_30_Days}} &= C_1 \cdot \log(1 + \text{frequency_doing_risky_things}) + C_2 \cdot \log(1 + \text{age}) \\U_{\text{Past_12_Months}} &= C_1 \cdot \text{painreliever_disorder_past_year} + C_1 \cdot \log(1 + \text{doctor_visits_last_year}) \\U_{\text{More_Than_12_Months}} &= C_1 \cdot \text{painreliever_disorder_past_year} + C_2 \cdot \text{self_evalute_mental_health} \\&\quad + C_3 \cdot \exp(-\text{education_level})\end{aligned}\quad (6)$$

3 Algorithm Details

Algorithm 1 Efficient Optimization with *DualMind Gradient*

Require: Training set $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$; group mapping $g(i) = G(\mathbf{x}_i)$; initial $\mathcal{H}^{(0)}, \{U_c^{(0)}\}$; LLMs f (decision), f' (updater), f_u (grad describer); templates τ, τ' ; epochs E

Ensure: Refined $\mathcal{H}^{(g)}$ and $\{U_c^{(g)}\}$

- 1: **for** $e = 1$ to E **do**
- 2: **for** each $(\mathbf{x}_i, y_i) \in \mathcal{D}$ **do**
- 3: $g \leftarrow G(\mathbf{x}_i)$
- 4: **// Forward: utilities and LLM predictions**
- 5: $T_i \leftarrow \tau(\mathbf{x}_i)$
- 6: $\hat{y}_i^{\text{base}} \leftarrow f(T_i)$ ▷ individual LLM prediction
- 7: **for** $c \in \mathcal{C}$ **do**
- 8: $u_{i,c}^{(g_t)} \leftarrow U_c^{(g_t)}(\mathbf{x}_i)$
- 9: **end for**
- 10: $\mathbf{u}_i^{(g_t)} \leftarrow (u_{i,1}^{(g_t)}, \dots, u_{i,C}^{(g_t)})$
- 11: $\boldsymbol{\pi}_i^{(g_t)} \leftarrow \text{SOFTMAX}(\mathbf{u}_i^{(g_t)})$
- 12: $\hat{y}_i^{\text{util}} \leftarrow \arg \max_{c \in \mathcal{C}} \boldsymbol{\pi}_{i,c}^{(g_t)}$
- 13: $\mathcal{I}_i^{(g_t)} \leftarrow \tau'(T_i, \mathcal{H}^{(g_t)}, \mathbf{u}_i^{(g_t)}, \boldsymbol{\pi}_i^{(g_t)})$
- 14: $\hat{y}_i^{\text{steer}} \leftarrow f(\mathcal{I}_i^{(g_t)})$ ▷ utility-steered LLM
- 15: **// Efficient DualMind Gradient Updates**
- 16: $c_{\text{LLM}} \leftarrow \mathbb{I}[\hat{y}_i^{\text{base}} = y_i]$
- 17: $c_{\text{util}} \leftarrow \mathbb{I}[\hat{y}_i^{\text{util}} = y_i]$
- 18: **if** $c_{\text{LLM}} = 1$ **and** $c_{\text{util}} = 1$ **then** ▷ both correct: skip
- 19: **continue**
- 20: **else if** $c_{\text{LLM}} = 0$ **and** $c_{\text{util}} = 1$ **then** ▷ update knowledge only
- 21: $\Delta \mathcal{H}^{(g_t)} \leftarrow \text{TEXTUALGRADIENT}(\mathcal{I}_i^{(g_t)}, y_i, \hat{y}_i^{\text{steer}}, \hat{y}_i^{\text{util}})$
- 22: $\mathcal{H}^{(g_{t+1})} \leftarrow f'(\mathcal{H}^{(g_t)}, \Delta \mathcal{H}^{(g_t)})$
- 23: **else if** $c_{\text{LLM}} = 1$ **and** $c_{\text{util}} = 0$ **then** ▷ update utilities only
- 24: $\Delta U^{(g_t)} \leftarrow \text{UTILITYGRADIENT}(\boldsymbol{\pi}_i^{(g_t)}, y_i, \{U_c^{(g_t)}\}, \mathbf{x}_i)$
- 25: $\text{grad_desc} \leftarrow f_u(\Delta U^{(g_t)})$
- 26: $\{U_c^{(g_{t+1})}\} \leftarrow f'(\{U_c^{(g_t)}\}, \text{grad_desc})$
- 27: **else** ▷ update both
- 28: $\Delta \mathcal{H}^{(g_t)} \leftarrow \text{TEXTUALGRADIENT}(\mathcal{I}_i^{(g_t)}, y_i, \hat{y}_i^{\text{steer}}, \hat{y}_i^{\text{util}})$
- 29: $\mathcal{H}^{(g_{t+1})} \leftarrow f'(\mathcal{H}^{(g_t)}, \Delta \mathcal{H}^{(g_t)})$
- 30: $\Delta U^{(g_t)} \leftarrow \text{UTILITYGRADIENT}(\boldsymbol{\pi}_i^{(g_t)}, y_i, \{U_c^{(g_t)}\}, \mathbf{x}_i)$
- 31: $\text{grad_desc} \leftarrow f_u(\Delta U^{(g_t)})$
- 32: $\{U_c^{(g_{t+1})}\} \leftarrow f'(\{U_c^{(g_t)}\}, \text{grad_desc})$
- 33: **end if**
- 34: **end for**
- 35: **end for**
- 36: **return** $\{\mathcal{H}^{(g)}, \{U_c^{(g)}\}\}_g$

4 Example of Co-Gradients

Here is the combined feedback we got for this specific Editable domain knowledge block, provide your feedback to this variable and other variables: To improve the prediction of <Self-employed><0.57> and address the feedback provided in the objective function, consider the following strategies for the editable utility function block and the editable domain knowledge block:

- **Expand Health and Economic Context:** The domain knowledge should include more comprehensive information about health and economic factors. For example, details about local job market conditions, availability of remote work opportunities, and economic support measures could provide a richer context for decision-making.

- **Incorporate Psychological and Social Factors:** Include knowledge about psychological factors such as stress levels, work-life balance preferences, and social influences. This could involve understanding how peer behavior or family expectations impact work choices. (...)

Supplementary Figure 1: Illustration of the co-gradient for the knowledge component in the JOB task.

Here is the combined feedback we got for this specific Editable utility function block, provide your feedback to this variable and other variables: To improve the prediction of <Self-employed><0.57> and address the feedback provided in the objective function, consider the following strategies for the editable utility function block and the editable domain knowledge block:

- **Incorporate Non-linear Relationships:** The current utility expressions are linear. Consider introducing non-linear terms or interaction effects to capture more complex relationships between variables. For example, the interaction between age and income might have a non-linear impact on the decision to be self-employed.

- **Dynamic Weighting:** Allow the weights in the utility function to be dynamic and context-dependent. For instance, the importance of 'covid_worked_from_home' might vary based on current pandemic conditions or personal health status. (...)

In addition to the textual gradient, please also consider the numeric gradient:

<Work for salary> - covid_worked_from_home: No clear effect on decision to work for salary (neutral effect).

Utility Form (Numeric): $1.95 * \text{covid_worked_from_home}$

Related Params: C_2=1.95 (Gradient: 0.00)

<Self-employed> - log(age + 1): Older individuals are more likely to choose self employment (positive effect).

Utility Form (Numeric): $0.97 * \log(\text{age} + 1)$

Related Params: C_2=0.97 (Gradient: +1.96) (...)

Supplementary Figure 2: Illustration of the co-gradient for the utility component in the JOB task.

5 Additional Results

5.1 Contributions of Grouped Features

To better understand how different components contribute to performance gains and interact during training, we decompose the input features across the six datasets into three categories: *cognitive* (belief-related factors), *context* (environmental factors), and *individual* (basic demographic attributes). We then quantify the contribution of each category throughout the optimization process. The settings are shown in Section 5.1.1 and the results are shown in Section 5.1.2.

5.1.1 Settings

Supplementary Table 3: Feature partition into Individual, Context, and Cognitive sets for different decision tasks.

Task	Individual	Context	Cognitive
Vaccine	age, gender, nurse, healthcare_worker, physician, have_university_degree, income_below_median, income_unknown	have_covid_sick_family_member, covid_threat, covid_preventable_by_vax, risk_of_covid_greater_than_vax	vaccine_safe_to_me, trust_government, trust_science, vax_protect_long_unsure, vax_protect_long_yes, less_attention_to_vax_info, more_attention_to_vax_info
Travel Mode	traveler_age_group, is_male, is_female, annual_income_level, is_first_class_traveler, ticket_payer_type, has_ga_travel_pass	trip_purpose, origin_canton_code, destination_canton_code, number_of_luggage_items, is_car_available	train_total_travel_time_min, train_ticket_cost_chf, train_service_headway_min, sm_travel_time_min, sm_ticket_cost_chf, sm_service_headway_min, car_travel_time_min, car_travel_cost_chf
Job	age, income, gender, self_rated_health, education, condition_diabetes, condition_heart_disease, condition_cancer, condition_obesity	covid_worked_from_home, restriction_childcare_closure, restriction_reduced_public_transport, restriction_work_from_home_req, restriction_quarantine_stayhome, tried_received_gov_assistance	mental_nervous_anxious_past7days, mental_depressed_past7days, mental_lonely_past7days, symptom_fever_past7days, symptom_sore_throat_past7days, symptom_headache_past7days, symptom_loss_taste_smell_past7days, covid_diagnosed_self, covid_diagnosed_household, swab_test_willing
Marijuana	age, education_level, income, gender, age_first_use, marijuana_disorder_past_year	state_marijuana_law, culty_getting_marijuana	diff- anxious_quitting_marijuana, sleep_issues_quitting_marijuana, strong_urge_marijuana, tried_to_stop_marijuana, risk_mj_weekly, opinion_adult_try_marijuana
Pain Reliever	age, education_level, income, gender, have_high_blood_pressure, have_asthma, num_children_in_household, painreliever_disorder_past_year	doctor_discuss_drug_use, doctor_visits_last_year, pital_overnight_last_year, times_attend_religious_services, religious_influence_decisions	doc- self_evaluate_mental_health, hos- ever_several_days_depressed, re- frequency_doing_risky_things
Alcohol	age, education_level, income, gender, have_high_blood_pressure, num_children_in_household, ever_binge_drinking, insurance_covers_alcohol_treatment	doctor_visits_last_year, place_drug_alcohol_policy, doctor_asked_how_much_alcohol	work- self_evaluate_mental_health, doc- frequency_doing_risky_things, ever_several_days_depressed, times_attend_religious_services, risk_harm_5plus_drinks_weekly

Prompt 5.1: System Prompt for Feature Categorization.

You are a feature categorization assistant.

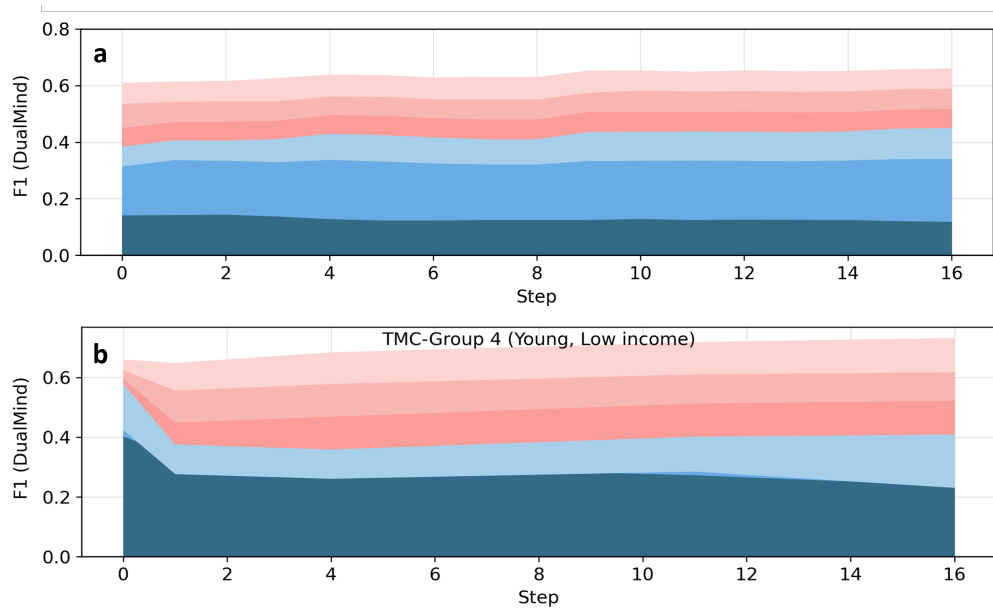
Task: You must categorize each bullet into one or more of the following three categories (multi-label allowed):

1. **Individual Attributes** – characteristics of the decision-maker, including demographic, socioeconomic, and personal health or resource factors. Focus on who the person is — background, capabilities, constraints.
2. **Context / Environment** – external or situational conditions surrounding the decision (geographic, policy, social, workplace or temporal environments) that shape options and pressures. Focus on where/under what circumstances behavior occurs.
3. **Cognitive / Attitudinal Factors** – mental, emotional, and perceptual aspects (beliefs, trust, motivation, perception of risk, psychological state). Focus on how/why people think, feel, and form intentions.

IMPORTANT:

- Multi-label is allowed and common. A bullet can belong to 1, 2, or all 3 categories.
- Return **STRICT JSON ONLY** with keys exactly: "individual", "context", "cognitive".
- Each value is an array of strings, where each string is one bullet **EXACTLY as provided** (no rewriting, no numbering).
- If a bullet fits multiple categories, include the exact same string in multiple arrays.

5.1.2 Results

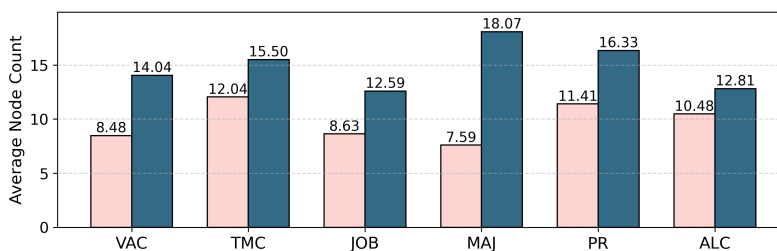


Supplementary Figure 3: Contributions of grouped features. (a) Aggregated feature contributions across datasets and optimization steps. We decompose input features into three categories: *cognitive* (belief-related factors), *context* (environmental factors), and *individual* (basic demographic attributes). Contributions of each component are computed at every step and aggregated to show overall trends. Results are averaged over the TMC, VAC, and JOB datasets, each optimized for 16 steps. (b) Evolution of the F1 score of the *DualMind* alongside contributions of the cognitive, context, and individual components for Group 4 (young, low-income) in the TMC dataset. Due to variability in API calls, the estimated multi-component contributions may differ slightly from the two-component (utility-LLM) values reported in the Results section.

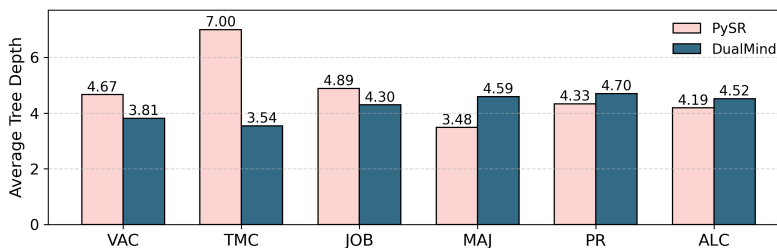
Utility drives the performance improvements. Figure 3a shows that, as optimization progresses, the F1 score increases from 0.610 to 0.660, driven primarily by strengthened utility-based components, while the contribution of the knowledge set remains relatively stable. In particular, the contextual and cognitive representational capacity of the utility-based components increases from 28.4% to 33.6% and from 11.6% to 16.6%, respectively. By aligning utility functions with textual knowledge, *DualMind* further unlocks the potential of utility-based representations, leading to consistent performance improvements driven by improved utility functions. Nevertheless, knowledge & concept set can drive the performance improvement for some group. In Figure 3b, we can observe the improvement is mainly brought by refined knowledge & concept set, where the contribution of all parts in the knowledge & concept set is increased.

5.2 Simplicity of Refined Utility Functions

Beyond predictive performance, *DualMind* also yields utility functions that are more interpretable. As illustrated in Figure 4 and 5, the refined utility functions exhibit a larger number of nodes while maintaining shallower tree depths, resulting in less nested and more transparent functional forms compared with PySR.



Supplementary Figure 4: Average node count of utility functions generated by PySR and *DualMind*.



Supplementary Figure 5: Average tree depth of utility functions generated by PySR and *DualMind*.

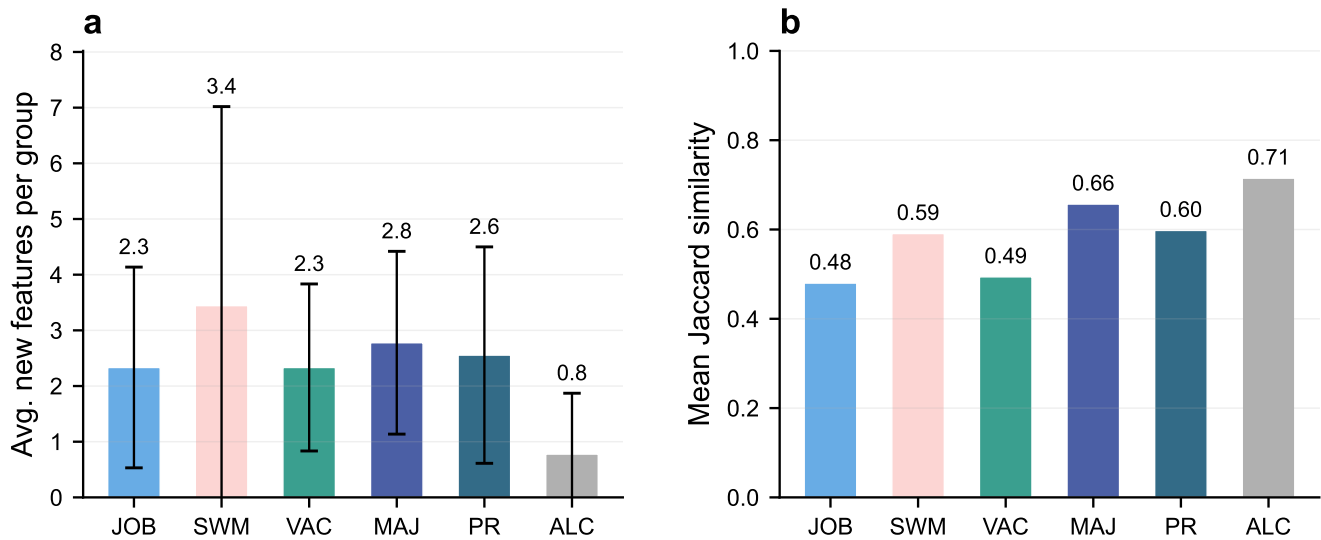
5.3 Feature Evolution in Utility Functions

Supplementary Table 4: Task-level summary of feature evolution. For each task, we report the number of newly discovered features across 9 demographic groups, alongside the mean pairwise Jaccard similarity of final feature sets. Higher Jaccard values indicate greater inter-group convergence.

Task	Total new	Avg/group	Std	Unique new	Mean Jaccard	Core (9/9)
JOB	21	2.3	1.7	12	0.48	3
SWM	31	3.4	3.4	8	0.59	3
VAC	21	2.3	1.4	8	0.49	2
MAJ	25	2.8	1.5	6	0.66	5
PR	23	2.6	1.8	8	0.60	4
ALC	7	0.8	1.0	5	0.71	4

Supplementary Table 5: Average number of new features by age group and task. Young groups consistently require more feature exploration than middle-aged and old groups, with the exception of SWM where old groups exhibit the highest demand for new travel-related features.

Age	JOB	SWM	VAC	MAJ	PR	ALC	Overall
Young	4.0	3.7	3.3	3.0	4.7	1.0	3.28
Middle	1.0	2.0	2.7	2.3	2.3	0.0	1.72
Old	2.0	4.7	1.0	3.0	0.7	1.3	2.11



Supplementary Figure 6: Task-level feature evolution characteristics. **a**, Average number of new features discovered per demographic group (\pm s.d.). SWM requires the most feature exploration (3.4 per group), while ALC is nearly solved by the initial specification (0.8 per group). **b**, Mean pairwise Jaccard similarity of final feature sets across 9 groups. ALC and MAJ show high inter-group convergence (Jaccard ≥ 0.66), whereas JOB and VAC exhibit substantial divergence (Jaccard ≤ 0.49), indicating that different demographic subgroups rely on distinct sets of features.

We use the Pairwise Jaccard Similarity (J) between the sets of selected utility features to quantify the similarity of learned utility feature sets across demographic groups. For two groups i and j , let F_i and F_j denote the sets of features selected in the learned utility functions. The Jaccard similarity is defined as:

$$J(F_i, F_j) = \frac{|F_i \cap F_j|}{|F_i \cup F_j|} \quad (7)$$

where $|\cdot|$ denotes the cardinality of a set. The pairwise Jaccard similarity ranges from 0 to 1, where higher values indicate greater overlap between the feature sets used by different groups. For each task, we report the mean pairwise Jaccard similarity across all demographic group pairs.

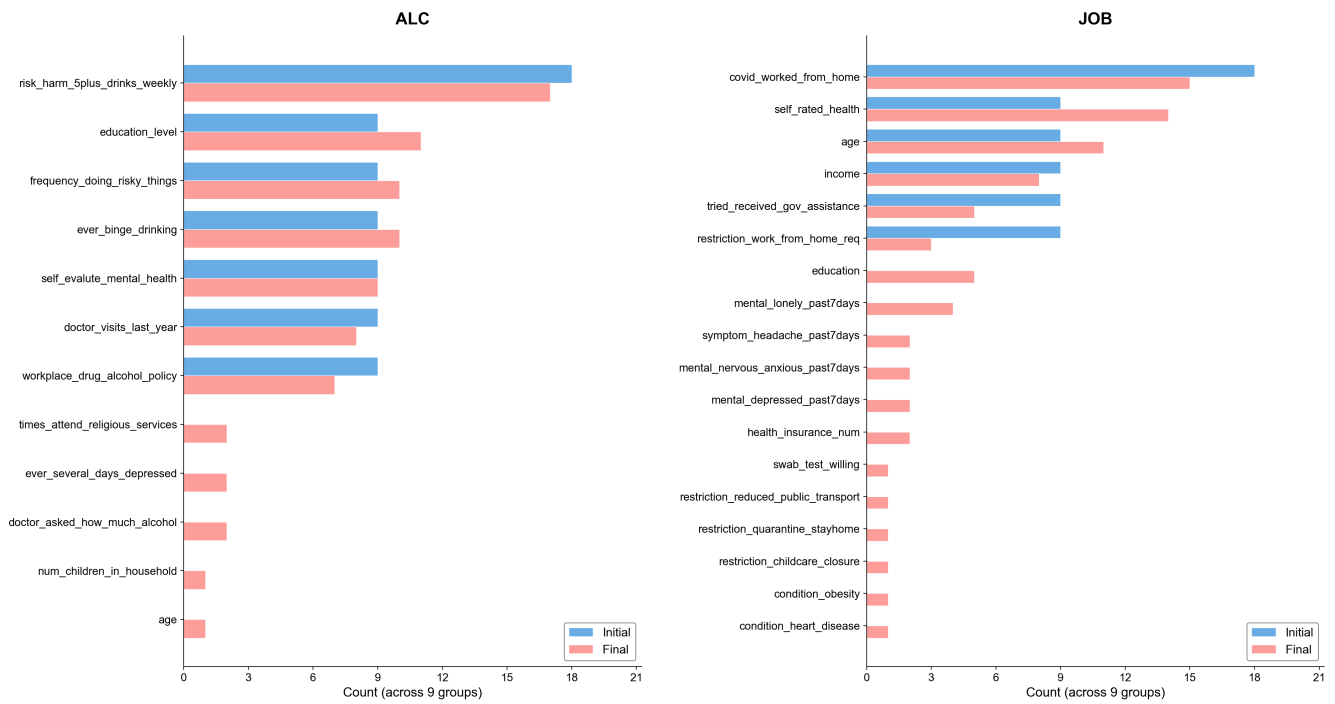
Feature discovery varies substantially across tasks. The optimization process introduces markedly different numbers of new features depending on the decision domain (Table 4, Figure 6a). SWM requires the most feature exploration, with an average of 3.4 new features per group—predominantly granular travel time and cost variables absent from the initial specification. In contrast, ALC needs only 0.8 new features per group on average, indicating that the initial feature set already captures the key determinants of drinking behavior. Tasks also differ in inter-group feature convergence (Figure 6b): ALC and MAJ exhibit high Jaccard similarity (≥ 0.66), meaning demographic



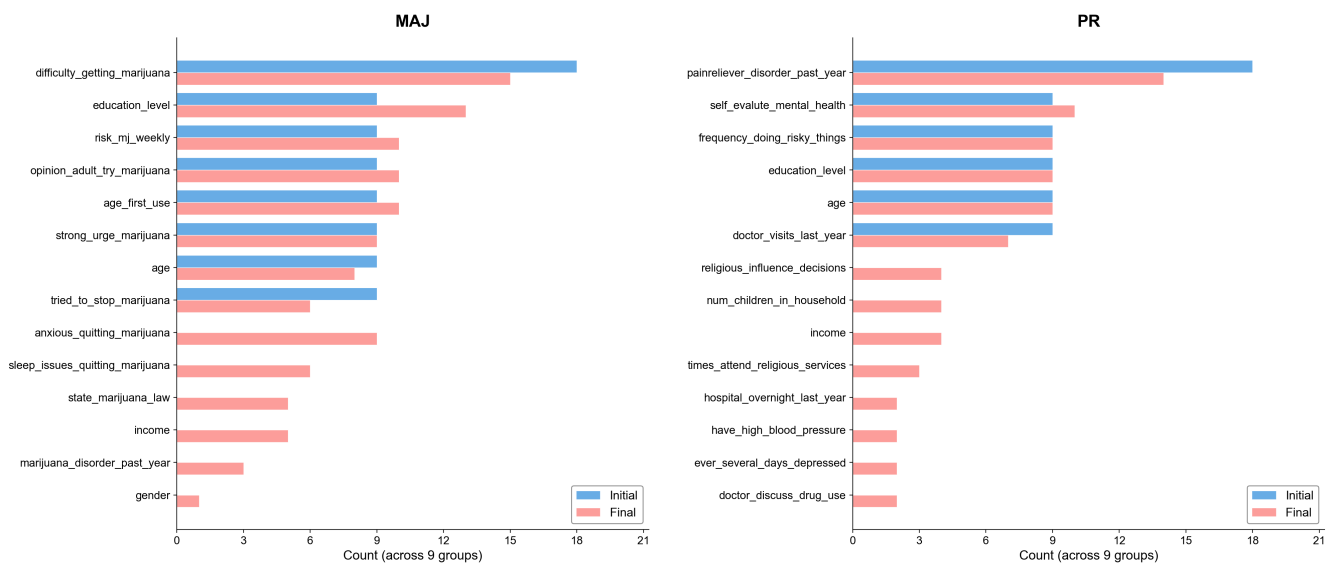
Supplementary Figure 7: Number of new features by demographic group and task. Each cell shows the count of features absent in the initial utility specification but present after optimization. Young groups (top block) consistently show higher counts across most tasks, particularly in PR (up to 6) and JOB (up to 5). SWM exhibits a reversed age pattern, with old groups requiring up to 8 new travel time/cost features. Gray cells indicate demographic–task combinations not present in the dataset.

subgroups converge on similar feature sets, while JOB and VAC show low similarity (≤ 0.49), suggesting that different populations rely on fundamentally different decision factors.

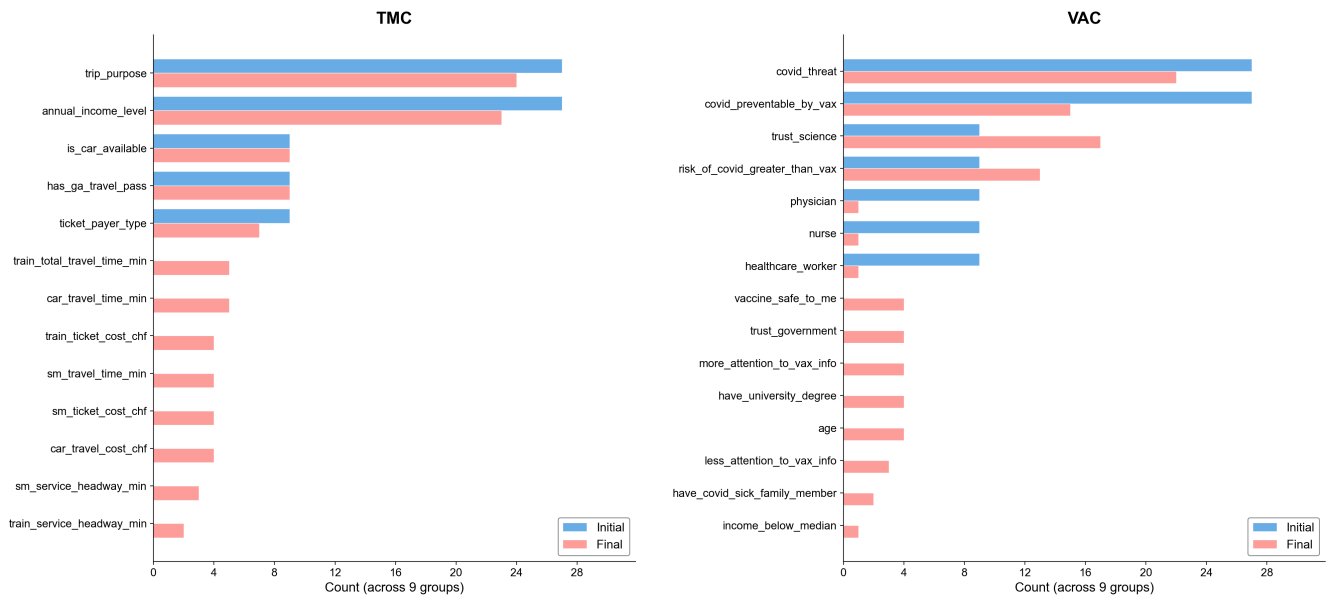
Young groups require more feature exploration. Across all six tasks, young demographic groups acquire substantially more new features than middle-aged groups (mean 3.28 vs. 1.72; Table 5, Figure 7). This pattern is especially pronounced in PR (pain reliever use), where young groups gain 4.7 new features on average compared to only 0.7 for old groups—suggesting that younger populations exhibit more complex decision-making processes that depend on a broader set of socioeconomic and contextual factors, whereas older populations rely on fewer, more stable determinants and are thus more readily captured by parsimonious utility specifications. A notable exception is SWM, where old groups require the most new features (4.7), reflecting the need for fine-grained travel cost and time attributes to model elderly mode choices.



Supplementary Figure 8: Distribution of utility terms in the initial and final states for ALC and JOB tasks.



Supplementary Figure 9: Distribution of utility terms in the initial and final states for MAJ and PR tasks.



Supplementary Figure 10: Distribution of utility terms in the initial and final states for TMC and VAC tasks.

5.4 Cost Analysis

5.4.1 Inference

The average inference cost per data sample using DualMind is reported in Table 6. We measure both the average token usage (T.N.) and the corresponding monetary cost (C) per sample for each dataset and LLM backbone. Across all settings, the per-sample inference cost of DualMind remains low and stable. For example, using GPT-4o-mini, DualMind requires approximately 1.3–1.5k tokens per sample, corresponding to a monetary cost of about $(0.21\text{--}0.23) \times 10^{-3}$ USD per inference. These results indicate that the proposed utility–LLM integration introduces only a small additional inference overhead while maintaining practical deployment efficiency.

Supplementary Table 6: Performance–efficiency comparison of Vanilla LLM, TextGrad, and DualMind during inference. TN and C denote the average token usage and monetary cost (in 10^{-3} USD) per LLM call for each data sample.

Methods	LLM Base	Vaccine			Travel Mode			Job			Marijuana			Pain Reliever			Alcohol		
		F1 ↑	TN ↓	C ↓	F1 ↑	TN ↓	C ↓	F1 ↑	TN ↓	C ↓	F1 ↑	TN ↓	C ↓	F1 ↑	TN ↓	C ↓	F1 ↑	TN ↓	C ↓
Vanilla LLM	GPT-4o-mini	0.354	1263	0.20	0.471	638	0.10	0.262	693	0.11	0.147	672	0.11	0.216	657	0.10	0.297	873	0.14
	GPT-5-mini	0.517	1273	0.21	0.506	642	0.10	0.383	701	0.11	0.330	680	0.11	0.431	668	0.11	0.265	880	0.14
	Gemini-2.5-FL	0.599	1293	0.40	0.414	682	0.07	0.263	721	0.07	0.231	703	0.07	0.000	690	0.08	0.235	912	0.10
	DeepSeek-V3.2	0.599	1288	0.36	0.445	670	0.19	0.341	705	0.20	0.194	688	0.20	0.422	679	0.19	0.230	894	0.25
TextGrad	GPT-4o-mini	0.480	1268	0.20	0.537	639	0.10	0.321	694	0.11	0.126	673	0.11	0.294	799	0.13	0.286	972	0.17
	GPT-5-mini	0.589	1273	0.21	0.508	643	0.10	0.223	697	0.11	0.219	681	0.11	0.445	758	0.13	0.335	1180	0.26
	Gemini-2.5-FL	0.522	1293	0.13	0.491	860	0.09	0.300	967	0.16	0.256	926	0.10	0.347	1526	0.30	0.346	1078	0.11
	DeepSeek-V3.2	0.480	1707	0.48	0.412	882	0.25	0.175	848	0.24	0.223	789	0.22	0.407	1048	0.30	0.203	1134	0.32
DualMind	GPT-4o-mini	0.662	1488	0.22	0.708	1480	0.23	0.664	1377	0.21	0.611	1449	0.22	0.552	1358	0.21	0.465	1368	0.21
	GPT-5-mini	0.666	1546	0.42	0.662	1444	0.39	0.635	1455	0.40	0.551	1453	0.41	0.549	1440	0.41	0.419	1396	0.40
	Gemini-2.5-FL	0.665	1556	0.16	0.653	1627	0.17	0.555	1508	0.15	0.586	1414	0.14	0.559	1328	0.14	0.422	1357	0.14
	DeepSeek-V3.2	0.671	1502	0.35	0.624	1538	0.36	0.626	1455	0.34	0.512	1410	0.33	0.559	1437	0.33	0.424	1406	0.33

5.4.2 Training

Supplementary Table 7: Training cost comparison between efficient and full optimization.

	VAC		MAJ	
	Efficient	Full	Efficient	Full
Time (wall-clock)	21m 29s	1h 9m 18s	23m 6s	50m 42s
API Cost (USD)	1.580	5.130	1.510	3.240
Input Tokens (4o-mini)	1,504,822	3,859,761	2,651,708	3,968,576
Input Tokens (4o)	285,559	972,070	239,880	569,730
Output Tokens (4o-mini)	6,794	17,063	25,771	37,714
Output Tokens (4o)	63,517	211,548	50,267	119,888
Optimization Steps	38	128	29	68
Mean F1	0.662	0.664	0.611	0.629

The training cost of *DualMind* under efficient and full update strategies is summarized in Table 7. We report wall-clock training time, token usage, API cost, and the number of optimization steps. The VAC and MAJ datasets are selected as representative cases to illustrate training cost, with 16 and 8 training iterations, respectively. Across both datasets, efficient updates substantially reduce training overhead compared with full updates. For example, on VAC, efficient updates require 21m 29s and \$1.58 of API cost, whereas full updates require 1h 9m 18s and \$5.13, corresponding to a $3.2\times$ reduction in time and $3.3\times$ reduction in cost. A similar trend holds for MAJ (23m 6s and \$1.51 vs. 50m 42s and \$3.24). This efficiency gain is consistent with the reduced optimization length (38 vs. 128 steps on VAC; 29 vs. 68 on MAJ). Despite the lower training cost, efficient updates achieve comparable performance to full updates (F1: 0.662 vs. 0.664 on VAC; 0.611 vs. 0.629 on MAJ).

5.5 Optimized Utility Functions and Knowledge Sets Examples

Utility Functions

<Unvaccinated>:

$0.574 \cdot \text{age} + 0.678 \cdot \text{covid_threat} - 0.396 \cdot \text{have_university_degree}$
 $+ 0.5 \cdot \text{income_below_median} + 0.113 \cdot \text{risk_of_covid_greater_than_vax}$
 $- 1.417 \cdot \log(\text{trust_science} + 1) + 0.031$

<Booster>:

$0.489 \cdot \text{age} + 0.28 \cdot \text{covid_threat} + 0.722 \cdot \text{have_university_degree}$
 $+ 0.668 \cdot \text{trust_science} + 0.549 \cdot \text{vaccine_safe_to_me} - 0.153$

<Vaccinate>:

$0.486 \cdot \text{age} + 1.221 \cdot \text{have_university_degree}$
 $+ 0.435 \cdot \text{risk_of_covid_greater_than_vax} + 0.601 \cdot \text{trust_science}$
 $+ 0.614 \cdot \text{vaccine_safe_to_me} + 0.106$

Knowledge Set

Vaccine Safety Belief

vaccine_safe_to_me significantly increases the likelihood of vaccination uptake.

Trust in Science & Government

trust_science, trust_government plays a critical role in vaccine uptake.

Risk Perception: COVID vs Vaccine

risk_of_covid_greater_than_vax enhances the decision to vaccinate.

Attention to Vaccine Information

more_attention_to_vax_info positively influences the decision to vaccinate.

Healthcare Worker Status

nurse, healthcare_worker, physician increases the likelihood of vaccination due to higher perceived risk and responsibility.

Supplementary Figure 11: Optimized utility functions and knowledge sets for VAC (young, high-income subgroup).

Utility Functions

<Unvaccinated>:

$0.383 * \text{covid_preventable_by_vax} + 1.215 * \text{covid_threat} +$
 $1.200 * \text{less_attention_to_vax_info} - 1.247 * \log(\text{trust_science} + 1) +$
 1.060

<Booster>:

$1.008 * \text{covid_preventable_by_vax} + 0.690 * \text{risk_of_covid_greater_than_vax}$
 $+ 0.825 * \text{covid_threat} - 0.441 * \text{have_covid_sick_family_member} - 0.095$

<Vaccinate>:

$0.924 * \text{covid_preventable_by_vax} + 0.979 * \text{covid_threat} +$
 $1.395 * \text{more_attention_to_vax_info} + 0.323 * \text{trust_government} - 0.315$

Knowledge Set

Risk & Threat Perception

Perceived COVID-19 threat and preventability by vaccination strongly increase vaccination preference.

Trust in Science & Government

Trust in scientific and governmental institutions promotes vaccination uptake.

Attention to Vaccine Information

Greater attention to vaccine information increases the likelihood of vaccination, whereas low attention is associated with remaining unvaccinated.

Perceived Relative Risk (COVID vs Vaccine)

Belief that COVID-19 risk exceeds vaccination risk encourages booster uptake.

Family COVID Experience

Having a family member sick with COVID-19 affects booster decisions.

Supplementary Figure 12: Optimized utility functions and associated knowledge sets for VAC (middle-, high-income subgroup).

Utility Functions

<Car>:

$1.614 * \text{annual_income_level} + 0.071 * \text{car_travel_cost_chf}$
 $- 0.029 * \text{car_travel_time_min} + 2.727 * \text{is_car_available} +$
 $1.299 * \text{trip_purpose} - 0.008$

<Train>:

$- 0.003 * \text{train_service_headway_min} + 0.082 * \text{train_ticket_cost_chf} -$
 $0.032 * \text{train_total_travel_time_min} + 1.663 * \text{trip_purpose} + 0.137$

<Swissmetro>:

$0.131 * \text{sm_service_headway_min} + 0.002 * \text{sm_ticket_cost_chf} +$
 $0.031 * \text{sm_travel_time_min} + 0.848 * \text{trip_purpose} + 0.001$

Knowledge Set

Car Availability

The availability of a car (`is_car_available`) decreases the likelihood of choosing Swissmetro.

Swissmetro Time & Cost

Lower Swissmetro travel time (`sm_travel_time_min`) and ticket cost (`sm_ticket_cost_chf`) increase the likelihood of choosing Swissmetro.

Service Frequency & Convenience

Frequent Swissmetro service (`sm_service_headway_min`) enhances convenience, increasing its attractiveness over other modes.

GA Travel Pass

Travelers with a GA travel pass (`has_ga_travel_pass`) are more likely to choose Swissmetro due to cost benefits.

Income Effects

High annual income (`annual_income_level`) may increase the likelihood of choosing first-class travel options, impacting Swissmetro choice.

Trip Purpose Heterogeneity

The purpose of the trip (`trip_purpose`) influences mode choice, with business travelers potentially favoring faster options like Swissmetro.

Regional Connectivity

The origin and destination canton codes (`origin_canton_code`, `destination_canton_code`) can affect travel mode choice based on regional connectivity and infrastructure.

Supplementary Figure 13: Optimized utility functions and associated knowledge sets for SWM (older, low-income subgroup).

Utility Functions

<Car>:

$0.838 \cdot \text{annual_income} - 0.003 \cdot \text{car_travel_time} + 2.268 \cdot \text{is_car_available} - 0.006 \cdot \text{trip_purpose} - 0.217$

<Train>:

$0.748 \cdot \text{annual_income} + 0.243 \cdot \text{ticket_payer_type} - 0.203 \cdot \text{trip_purpose} + 0.085 \cdot \sqrt{\text{train_ticket_cost}} + 0.108$

<Swissmetro>:

$1.613 \cdot \text{annual_income} + 0.054 \cdot \text{sm_service_headway} + 0.739 \cdot \sin(\text{trip_purpose}) + 0.233$

Knowledge Set

Car Availability

Car availability significantly increases the likelihood of choosing a car as the travel mode.

Travel Time

Travel time is a critical factor in determining the preference for car travel.

Income Preference

Higher annual income levels can lead to a preference for Swissmetro due to its premium service offerings.

Service Headway

Train and Swissmetro service headway influences the convenience and attractiveness of rail travel.

Ticket Cost

Train ticket cost directly affects the attractiveness of train travel.

Luggage

The number of luggage items influences mode choice, with fewer items favoring Swissmetro for ease of travel.

Supplementary Figure 14: Optimized utility functions and associated knowledge sets for SWM (young, low-income subgroup).

Utility Functions

<Work for salary>:

$2.142 * \text{covid_worked_from_home} + 0.731 * \text{restriction_work_from_home_req} - 0.046 * \log(\text{income} + 1) - 0.593$

<Self-employed>:

$0.606 * \text{covid_worked_from_home} - 0.925 * \text{symptom_headache_past7days} + 0.418 * \log(\text{age} + 1) - 0.363$

<Did not work>:

$- 0.614 * \text{covid_worked_from_home} + 0.247 * \text{self_rated_health} - 0.355 * \text{tried_received_gov_assistance} + 0.816$

Knowledge Set

COVID Health Status

Personal or family COVID diagnosis and related health conditions influence the decision to work for salary by affecting health status and availability.

Health and Mental Symptoms

Self-rated health and recent symptoms, such as headache or other mental and physical health conditions, can reduce the likelihood of working for salary.

Income Buffering Effect

Higher income levels may mitigate health-related work constraints, increasing the likelihood of working for salary.

Supplementary Figure 15: Optimized utility functions and associated knowledge sets for JOB (middle, high-income subgroup).

Utility Functions

<Work for salary>:

$1.636 * \text{covid_worked_from_home} + 0.360 * \text{self_rated_health} + 0.316 * \log(\text{income} + 1) - 0.082$

<Self-employed>:

$0.230 * \text{mental_lonely_past7days} + 0.710 * \text{self_rated_health} + 0.420 * \log(\text{age} + 1) - 0.265$

<Did not work>:

$0.467 * \text{condition_heart_disease} + 0.634 * \text{self_rated_health} + 0.402 * \log(\text{age} + 1) + 0.189$

Knowledge Set

Education and Remote Work

Higher education increases the likelihood of remote work opportunities.

Age and Work Adaptation

Age influences work behavior, with younger individuals more likely to adapt to remote work.

Health and Work Capacity

Self-rated health impacts work decisions, with better health increasing the likelihood of working from home.

Mental Health Constraints

Mental health factors, such as anxiety and depression, can decrease the likelihood of choosing remote work.

COVID Diagnosis Effects

Covid-19 diagnosis, either self or household, increases the likelihood of remote work or not working.

Policy Restrictions

Restrictions like work-from-home requirements directly increase remote work likelihood.

Income and Work Options

Income level affects work choices, with higher income individuals more likely to have remote work options.

Supplementary Figure 16: Optimized utility functions and knowledge sets for JOB (old, high-income subgroup).

Utility Functions

<Past_30_Days>:

$0.428 * \sqrt{\text{difficulty_getting_marijuana}} + 2.528 * \text{strong_urge_marijuana} + 0.385 * \log(\text{age_first_use} + 1) - 0.621$

<Past_12_Months>:

$0.438 * \text{anxious_quitting_marijuana} + 0.407 * \text{difficulty_getting_marijuana} + 1.388 * \exp(-\text{education_level}) - 0.567$

<More_Than_12_Months>:

$- 0.423 * \text{opinion_adult_try_marijuana} - 0.212 * \text{risk_mj_weekly} + 1.374 * \text{state_marijuana_law} + 0.605$

Knowledge Set

Age of First Use

Earlier age of first marijuana use is associated with shorter intervals between use occurrences.

Policy Environment

State marijuana laws can influence recency intervals by affecting accessibility and social acceptance.

Urge Intensity

Strong urges to use marijuana increase the likelihood of shorter recency intervals.

Access Difficulty

Difficulty in obtaining marijuana can lead to longer intervals between use.

Social Acceptance

Opinions on adult marijuana use can affect social acceptance, influencing recency intervals.

Supplementary Figure 17: Optimized utility functions and knowledge sets for MAJ (young, high-income subgroup).

Utility Functions

<Past_30_Days>:

$-0.175 * \text{anxious_quitting_marijuana} + 0.384 * \text{difficulty_getting_marijuana}$
 $+ 0.830 * \text{sleep_issues_quitting_marijuana} + 1.194 * \text{strong_urge_marijuana}$
 $+ 0.652 * \log(\text{age_first_use} + 1) - 0.253$

<Past_12_Months>:

$0.603 * \text{difficulty_getting_marijuana} + 0.118 * \text{marijuana_disorder_past_year}$
 $+ 1.120 * \text{tried_to_stop_marijuana} + 2.465 * \exp(-\text{education_level}) + 0.400$

<More Than 12 Months>:

$0.400 * \text{opinion_adult_try_marijuana} + 1.492 * \text{risk_mj_weekly} +$
 $0.205 * \text{state_marijuana_law} + 0.362 * \cos(\text{age}) + 0.086$

Knowledge Set

Age of First Use

Earlier initiation of marijuana use (age_first_use) can be associated with a higher likelihood of frequent use.

Policy Environment

State marijuana laws (state_marijuana_law) can influence the frequency and recency of marijuana use.

Urge Intensity

Strong urges to use marijuana (strong_urge_marijuana) increase the likelihood of frequent use.

Access Difficulty

Difficulty in obtaining marijuana (difficulty_getting_marijuana) can decrease the likelihood of frequent use.

Social Acceptance

Opinions on adult marijuana use (opinion_adult_try_marijuana) can shape individual usage patterns.

Withdrawal Anxiety

Anxious feelings when quitting marijuana (anxious_quitting_marijuana) may lead to increased usage frequency.

Withdrawal Sleep Issues

Sleep issues when quitting marijuana (sleep_issues_quitting_marijuana) can contribute to continued use.

Cessation Attempts

Attempts to stop marijuana use (tried_to_stop_marijuana) may indicate higher usage frequency.

Risk Perception

Risk perception of weekly marijuana use (risk_mj_weekly) can influence usage patterns.

Utility Functions

<Past_30_Days>:

$0.037 \cdot \log(\text{age} + 1) + 0.151 \cdot \log(\text{frequency_doing_risky_things} + 1) + 0.277 \cdot \log(\text{doctor_visits_last_year} + 1) - 0.985 \cdot \text{have_high_blood_pressure} - 0.216$

<Past_12_Months>:

$0.850 \cdot \text{painreliever_disorder_past_year} + 1.035 \cdot \exp(-\text{education_level}) - 0.835 \cdot \log(\text{num_children_in_household} + 1) + 0.305$

<More_Than_12_Months>:

$0.053 \cdot \text{religious_influence_decisions} + 0.721 \cdot \text{self_evalute_mental_health} + 0.084 \cdot \log(\text{income} + 1) - 0.203$

Knowledge Set

Life Events

Recent life events, such as changes in employment status, can influence the patterns and frequency of pain reliever usage.

Behavioral Health Indicators

Sleep quality and exercise frequency are behavioral indicators that can impact an individual's propensity to misuse pain relievers.

Contextual Influences

Seasonal changes and social influences are contextual variables that may affect pain reliever usage patterns.

Supplementary Figure 19: Optimized utility functions and knowledge sets for PR (young, low-income subgroup).

Utility Functions

<Past_30_Days>:

$0.374 \cdot \log(\text{age} + 1) - 0.437 \cdot \log(\text{doctor_visits_last_year} + 1) - 0.118 \cdot \log(\text{frequency_doing_risky_things} + 1) + 0.143$

<Past_12_Months>:

$0.913 \cdot \text{painreliever_disorder_past_year} + 0.117 \cdot \log(\text{income} + 1) - 0.620 \cdot \exp(-\text{education_level}) + 0.316$

<More_Than_12_Months>:

$1.199 \cdot \text{self_evalute_mental_health} + 0.456 \cdot \log(\text{num_children_in_household} + 1) + 0.040 \cdot \log(\text{times_attend_religious_services} + 1) - 0.404$

Knowledge Set

Temporal Patterns

Temporal patterns, such as seasonal variations, can influence pain reliever use, with certain times of the year potentially increasing usage.

Behavioral Factors

Behavioral insights, including medication adherence and psychological barriers, can affect the likelihood of pain reliever misuse.

Medical Context

Contextual factors like recent doctor visits or hospital stays may increase the likelihood of pain reliever use.

Mental Health

Self-evaluation of mental health can impact the frequency of pain reliever use, with poorer self-evaluations potentially leading to higher usage.

Religious Influence

Religious influence on decisions may decrease the likelihood of engaging in risky behaviors, including misuse of pain relievers.

Supplementary Figure 20: Optimized utility functions and knowledge sets for PR (middle, low-income subgroup).

Utility Functions

<Past_30_Days>:

$2.773 * \text{doctor_asked_how_much_alcohol} + 1.842 * \text{ever_binge_drinking} + 0.165 * \text{ever_several_days_depressed} + 0.401 * \text{frequency_doing_risky_things} - 1.395 * \log(\text{doctor_visits_last_year} + 1) + 0.039$

<Past_12_Months>:

$0.213 * \text{ever_binge_drinking} + 0.092 * \text{ever_several_days_depressed} - 0.056 * \text{frequency_doing_risky_things} + 0.408 * \text{risk_harm_5plus_drinks_weekly} - 0.063$

<More_Than_12_Months>:

$0.546 * \text{num_children_in_household} + 0.408 * \text{risk_harm_5plus_drinks_weekly} + 3.626 * \exp(-\text{education_level}) + 0.083$

Knowledge Set

Risky Behavior

Higher frequency of risky behaviors increases the likelihood of more recent alcohol consumption.

Mental Health

Poor self-evaluation of mental health is associated with more recent drinking behavior.

Medical Monitoring

Regular doctor visits may decrease the likelihood of recent alcohol misuse.

Religious Attendance

Attendance at religious services is often linked to less recent alcohol consumption.

Harm Perception

Perception of harm from frequent drinking can reduce the recency of alcohol use.

Supplementary Figure 21: Optimized utility functions and knowledge sets for ALC (young, high-income subgroup).

Utility Functions

<Past_30_Days>:

$0.875 * \text{ever_binge_drinking} + 1.215 * \text{frequency_doing_risky_things} + 1.633 * \text{self_evalute_mental_health} - 0.482$

<Past_12_Months>:

$0.319 * \text{risk_harm_5plus_drinks_weekly} - 0.695 * \text{self_evalute_mental_health} - 0.634 * \text{workplace_drug_alcohol_policy} + 0.048 * \exp(-\text{education_level}) + 0.334$

<More_Than_12_Months>:

$0.214 * \text{risk_harm_5plus_drinks_weekly} + 1.590 * \text{self_evalute_mental_health} + 0.166 * \log(\text{age}) - 0.134$

Knowledge Set

Medical Contact

Higher frequency of doctor visits in the last year is often associated with shorter intervals since the last alcohol use.

Mental Health

Self-evaluation of poor mental health increases the likelihood of more recent alcohol use.

Religious Attendance

Attendance at religious services is often associated with longer intervals since the last alcohol use.

Workplace Policy

Having a workplace drug and alcohol policy is often associated with longer intervals since the last alcohol use.

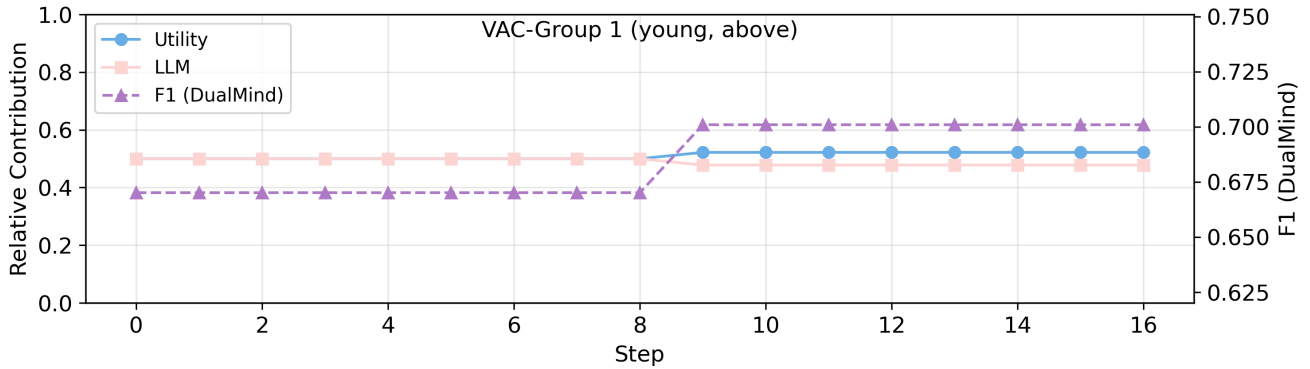
Treatment Coverage

Insurance coverage for alcohol treatment increases the likelihood of longer intervals since the last alcohol use.

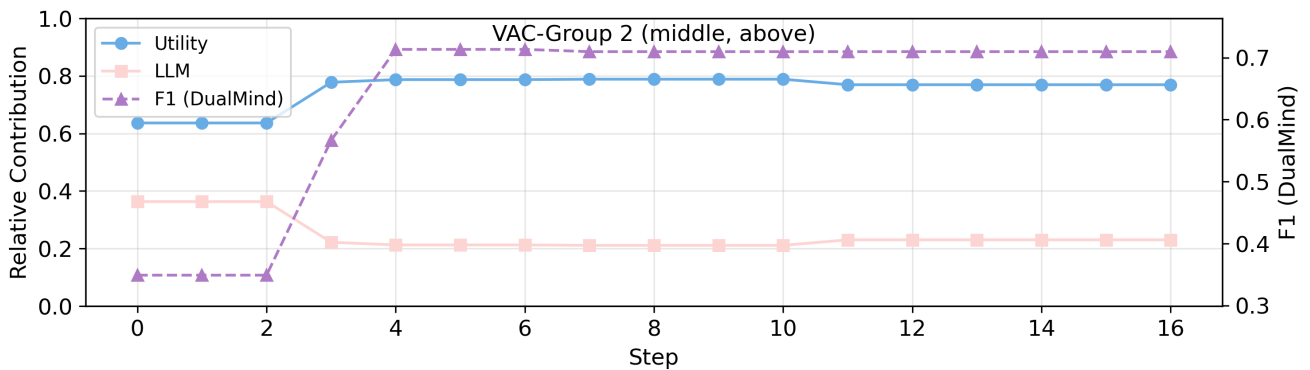
Supplementary Figure 22: Optimized utility functions and knowledge sets for ALC (old, middle-income subgroup).

5.6 Utility-Knowledge Contributions for All Groups and Tasks

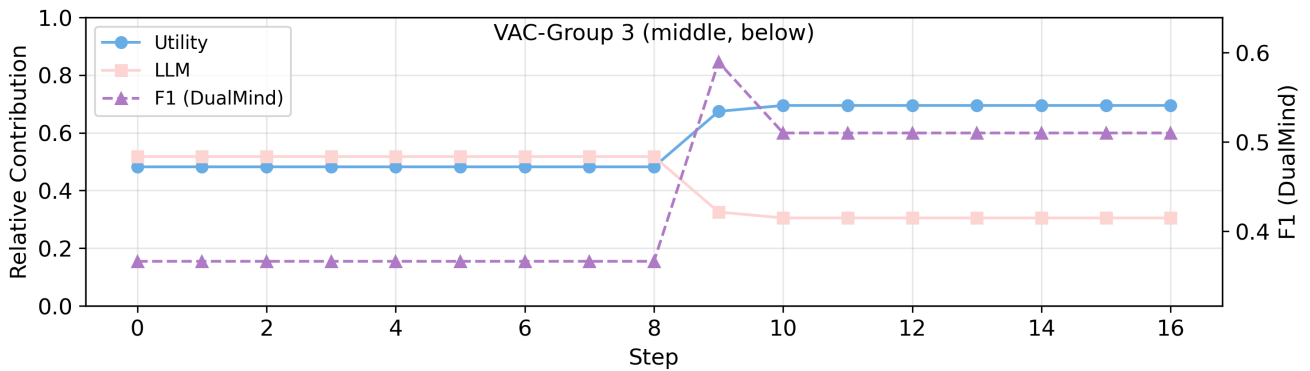
5.6.1 VAC



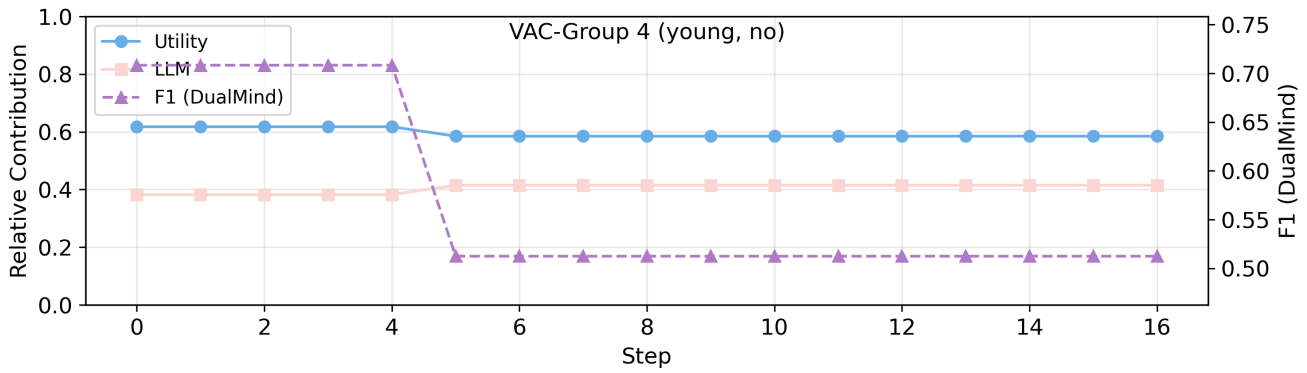
Supplementary Figure 23: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for VAC (young, High-income) in the VAC dataset.



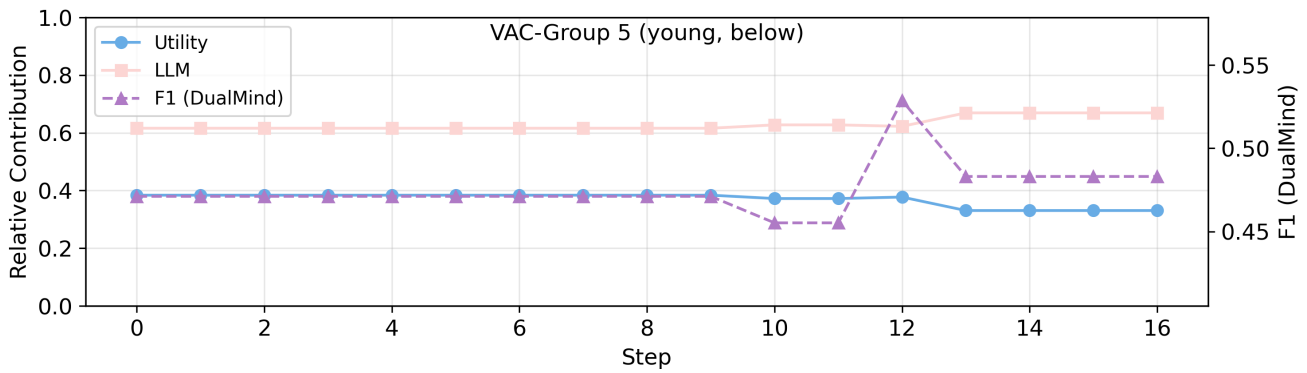
Supplementary Figure 24: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for VAC (middle, High-income) in the VAC dataset.



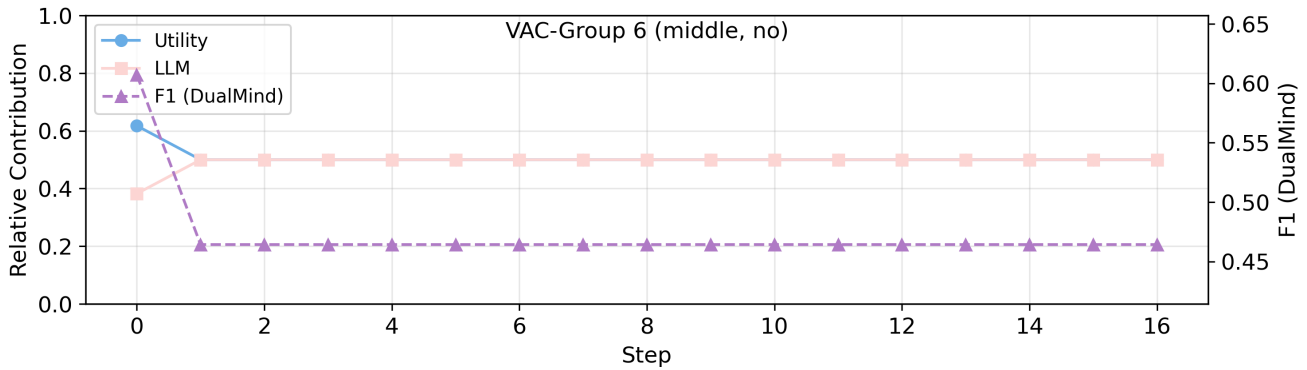
Supplementary Figure 25: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for VAC (middle, Low-income) in the VAC dataset.



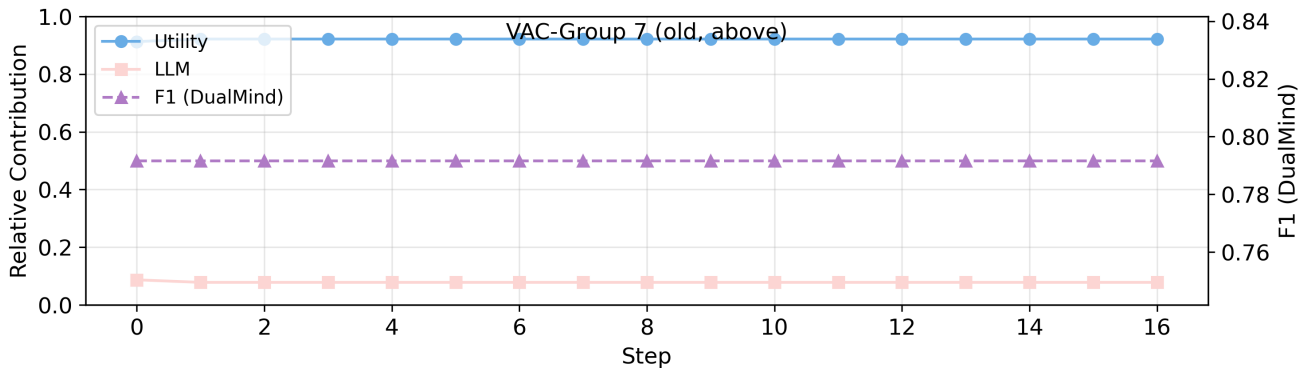
Supplementary Figure 26: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for VAC (young, Not specified) in the VAC dataset.



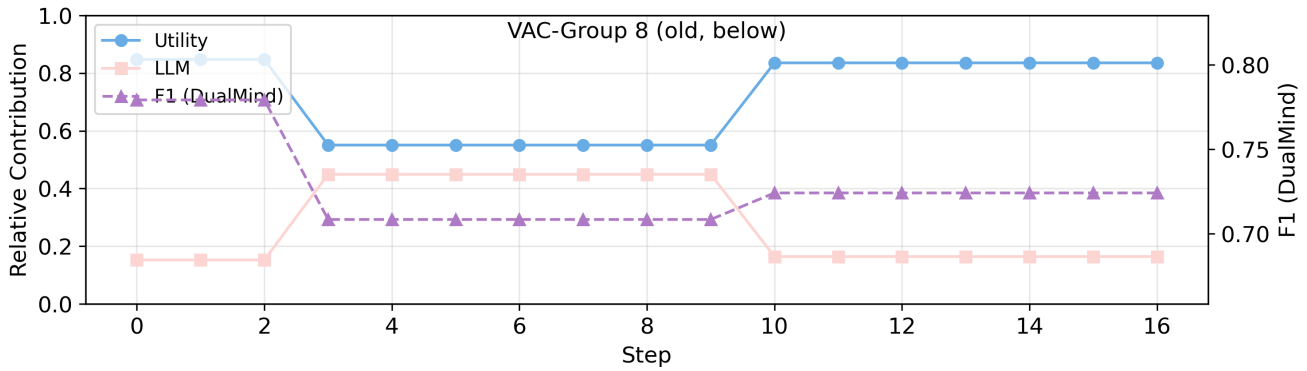
Supplementary Figure 27: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for VAC (young, Low-income) in the VAC dataset.



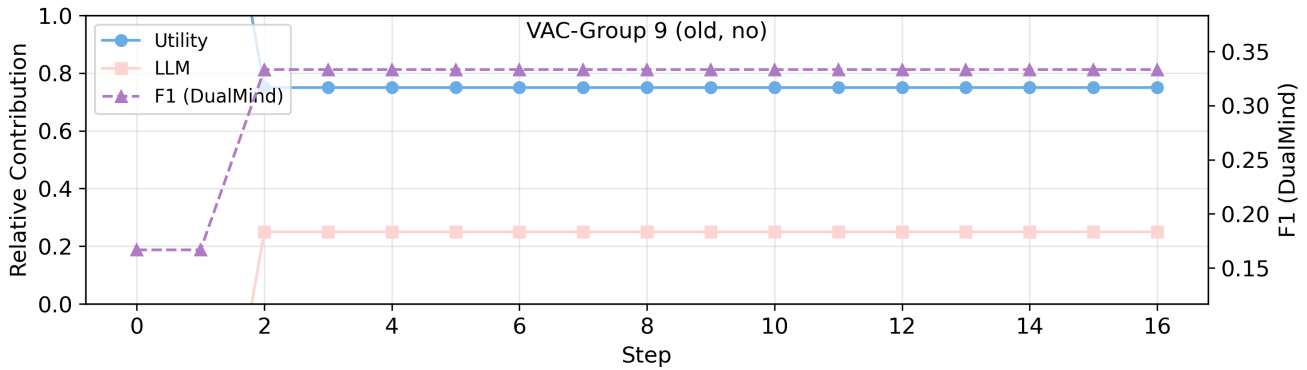
Supplementary Figure 28: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for VAC (middle, Not specified) in the VAC dataset.



Supplementary Figure 29: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for VAC (old, High-income) in the VAC dataset.

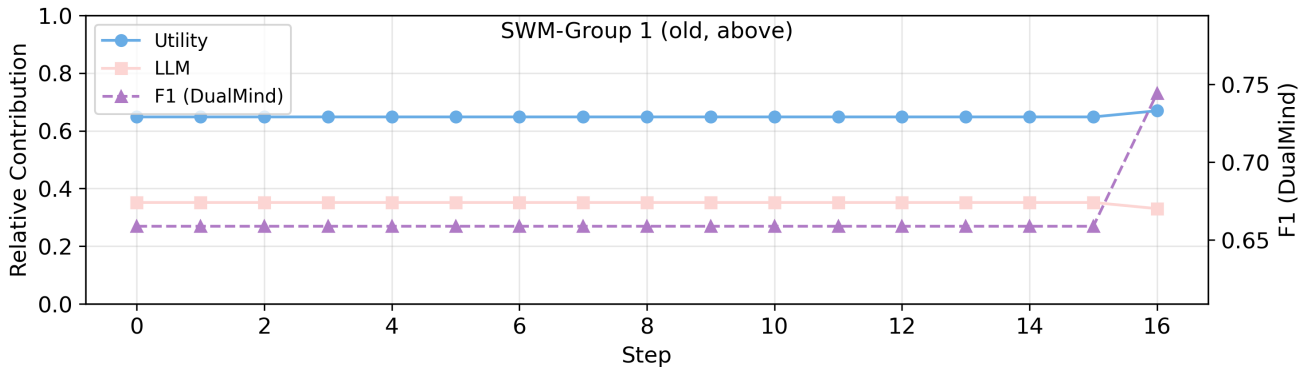


Supplementary Figure 30: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for VAC (old, Low-income) in the VAC dataset.

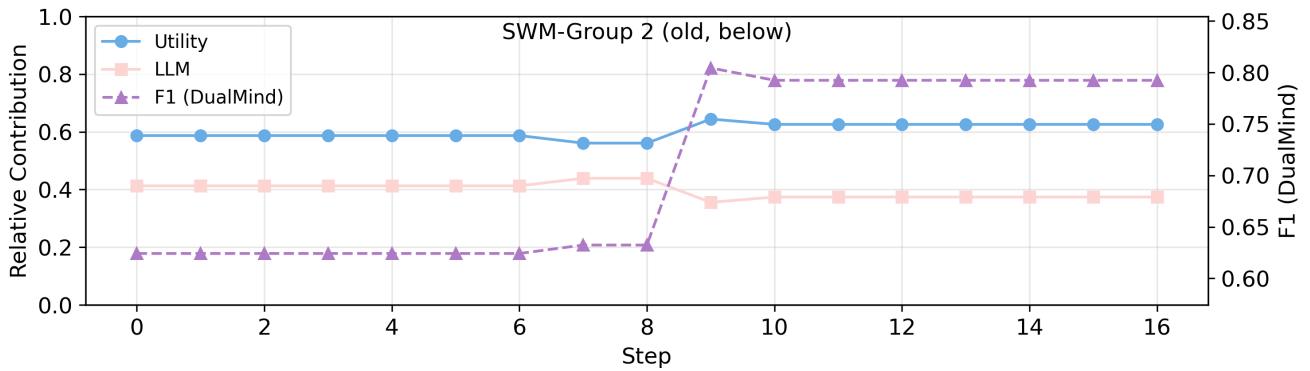


Supplementary Figure 31: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for VAC (old, Not specified) in the VAC dataset.

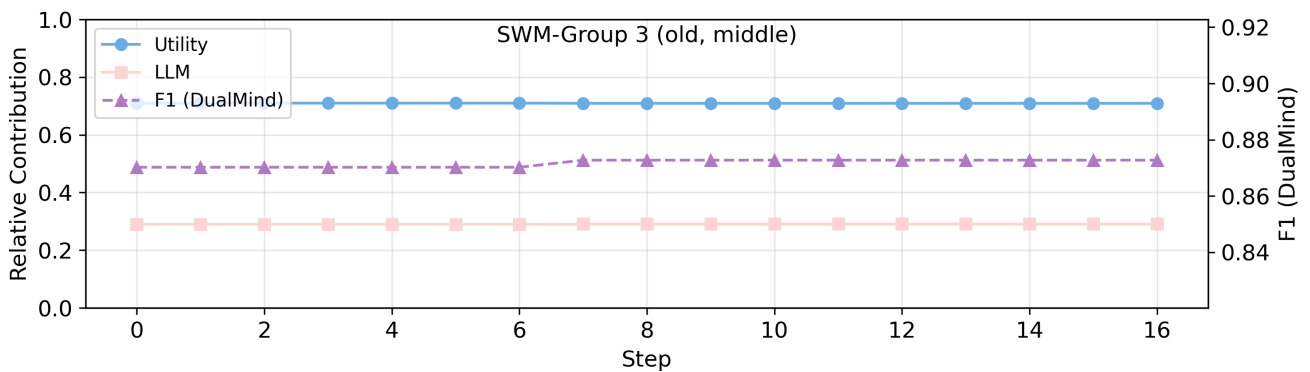
5.6.2 TMC



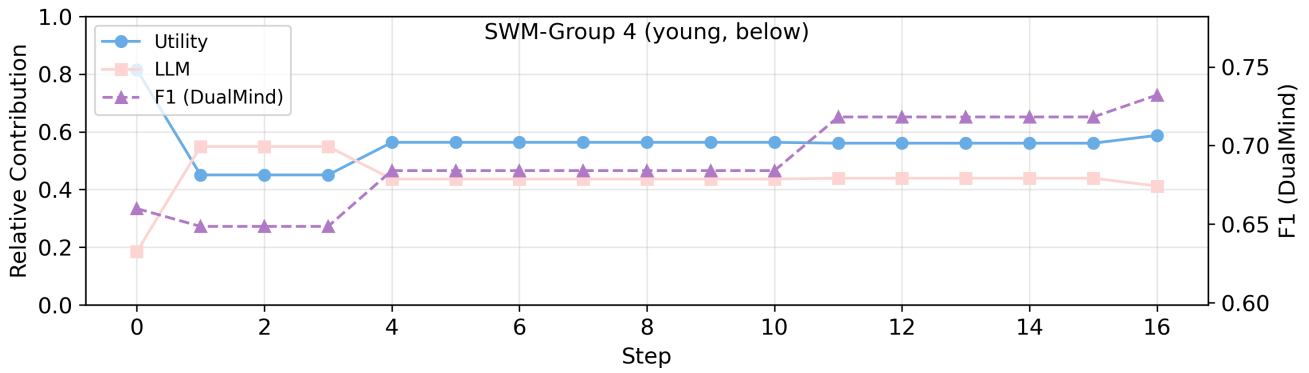
Supplementary Figure 32: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for TMC (old, High-income) in the TMC dataset.



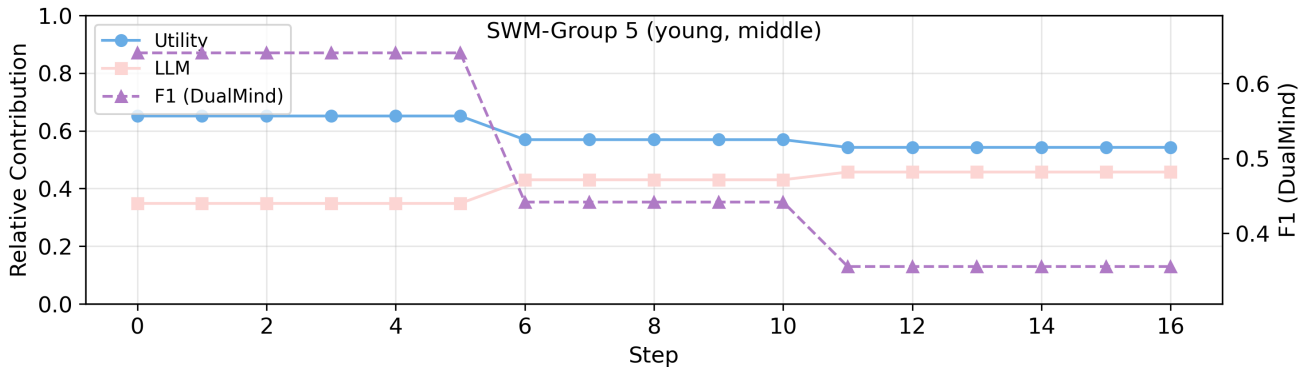
Supplementary Figure 33: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for TMC (old, Low-income) in the TMC dataset.



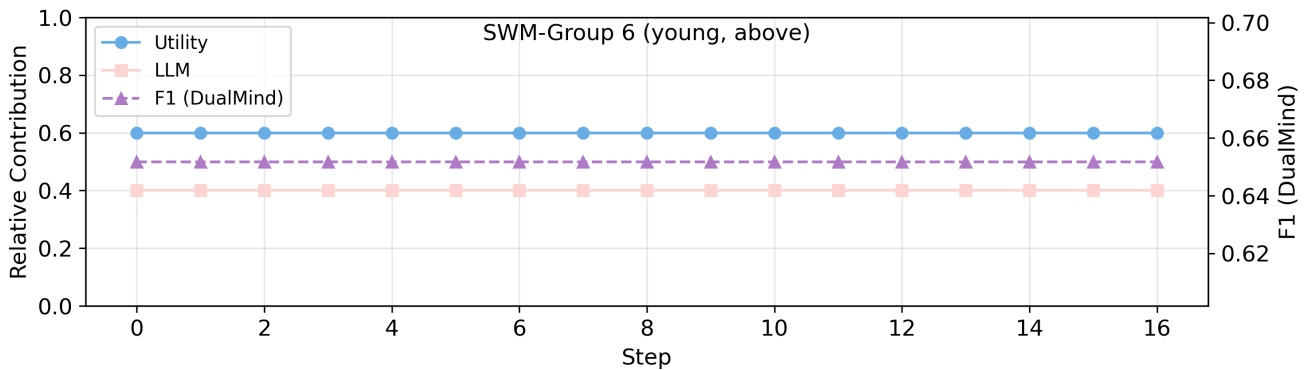
Supplementary Figure 34: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for TMC (old, Middle-income) in the TMC dataset.



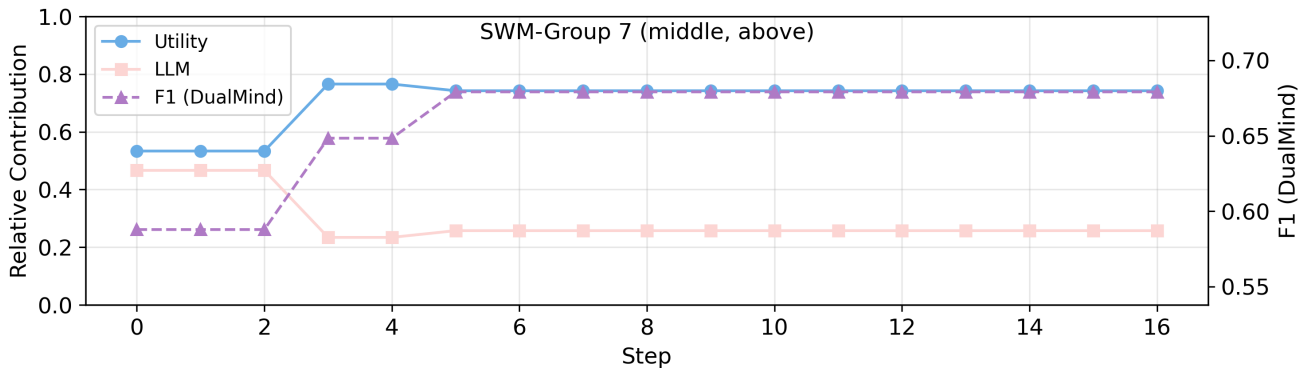
Supplementary Figure 35: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for TMC (young, Low-income) in the TMC dataset.



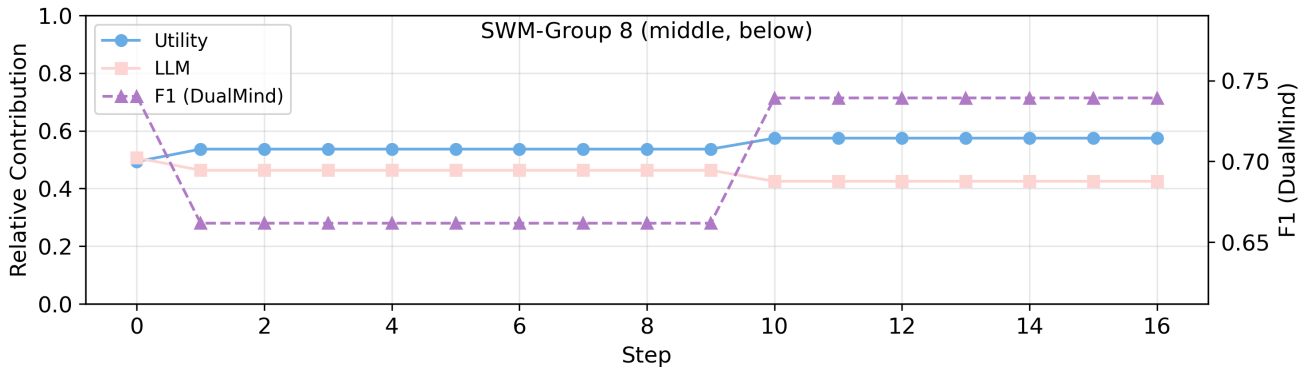
Supplementary Figure 36: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for TMC (young, Middle-income) in the TMC dataset.



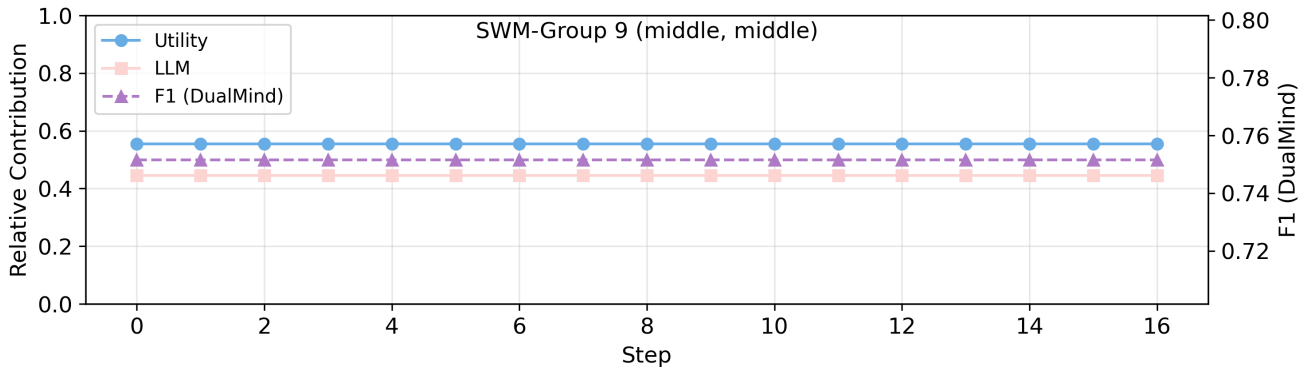
Supplementary Figure 37: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for TMC (young, High-income) in the TMC dataset.



Supplementary Figure 38: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for TMC (middle, High-income) in the TMC dataset.

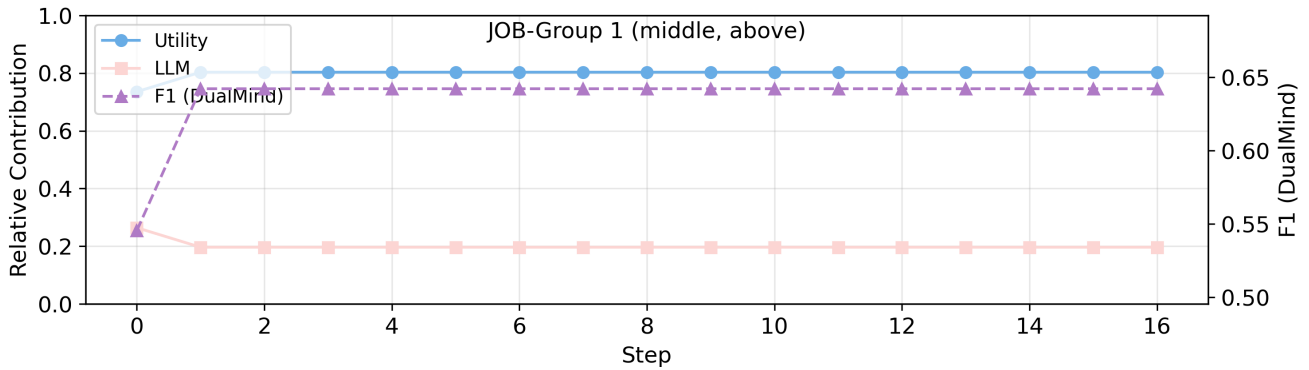


Supplementary Figure 39: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for TMC (middle, Low-income) in the TMC dataset.

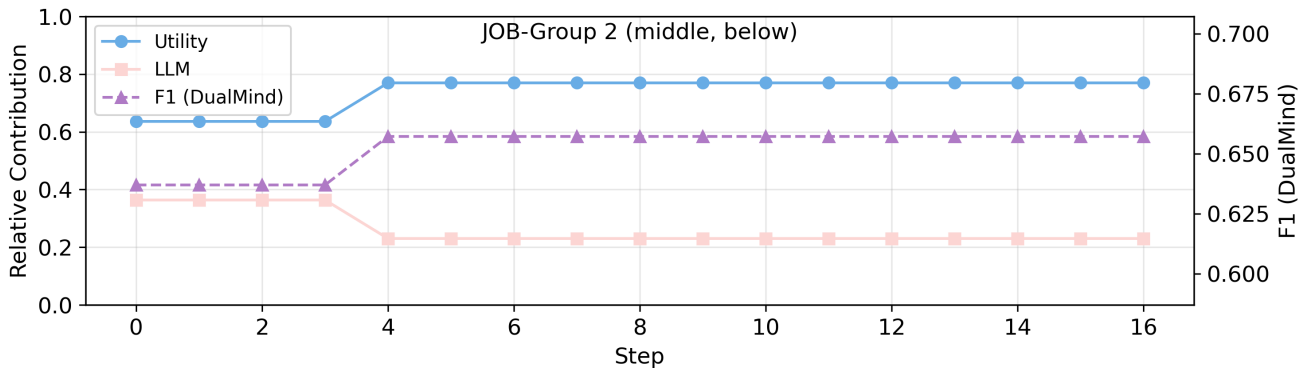


Supplementary Figure 40: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for TMC (middle, Middle-income) in the TMC dataset.

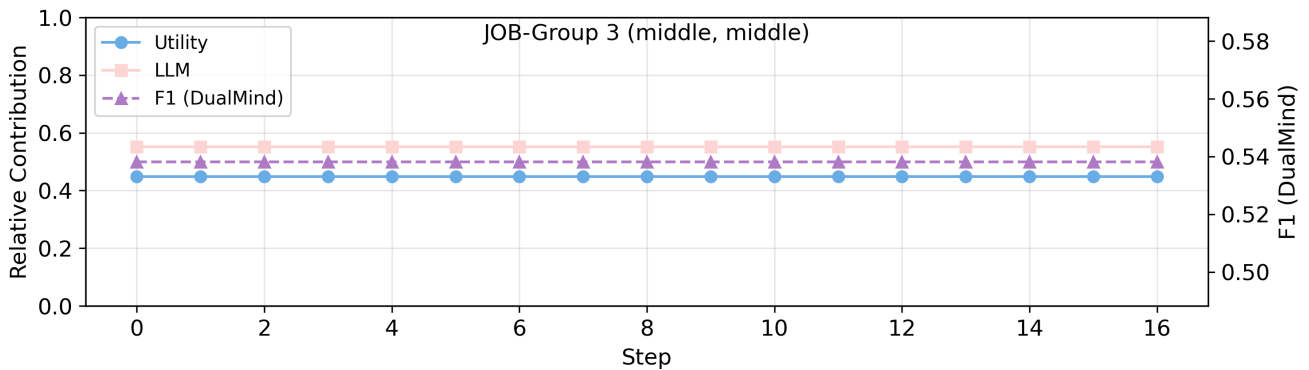
5.6.3 JOB



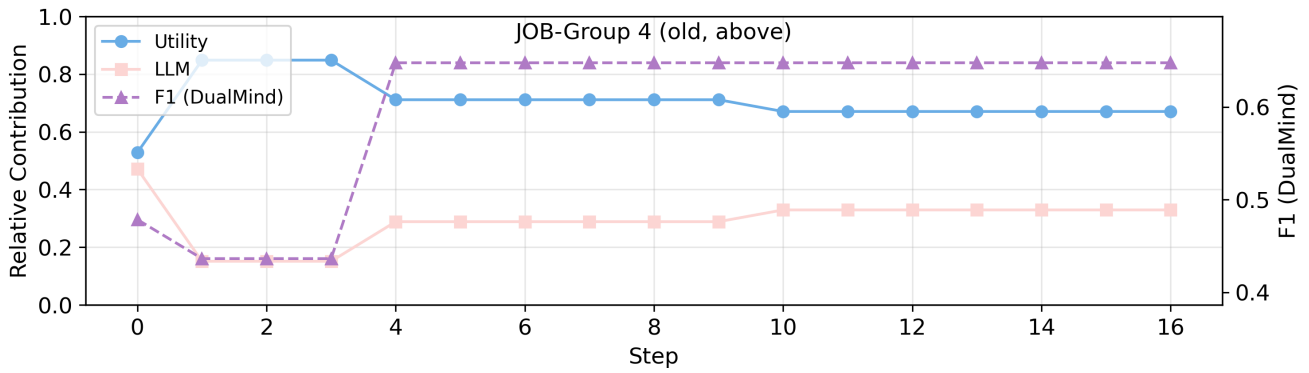
Supplementary Figure 41: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for JOB (middle, High-income) in the JOB dataset.



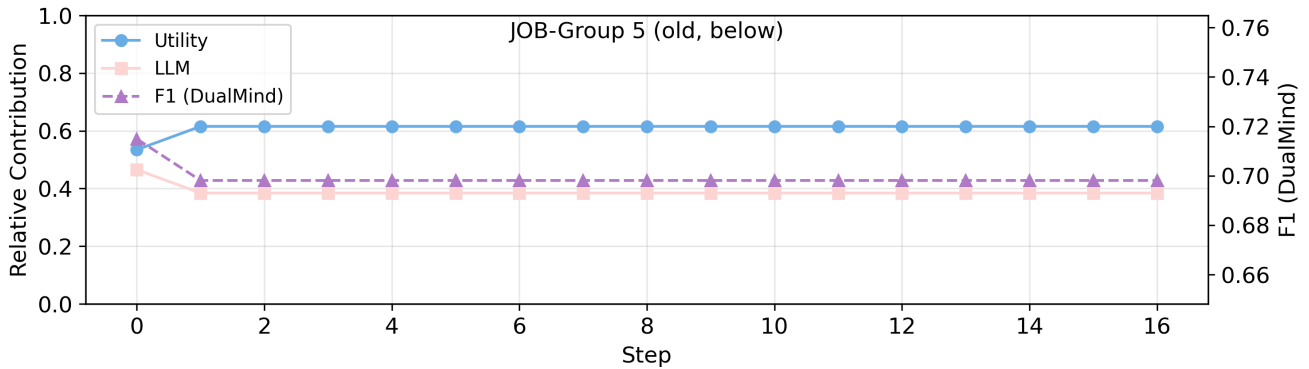
Supplementary Figure 42: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for JOB (middle, Low-income) in the JOB dataset.



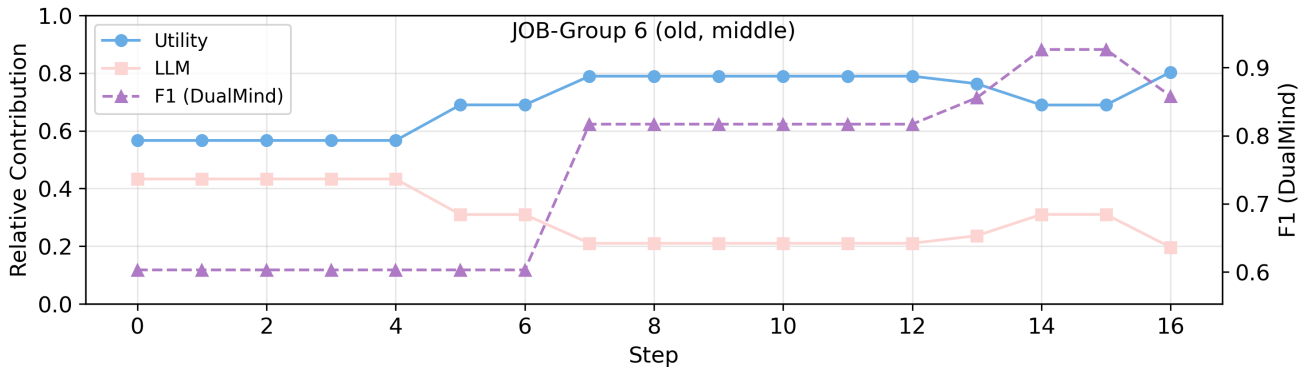
Supplementary Figure 43: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for JOB (middle, Middle-income) in the JOB dataset.



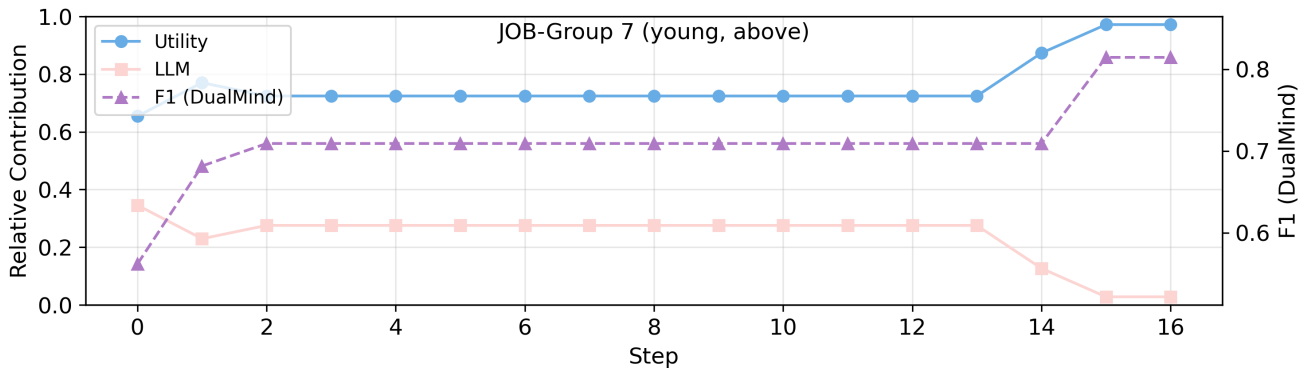
Supplementary Figure 44: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for JOB (old, High-income) in the JOB dataset.



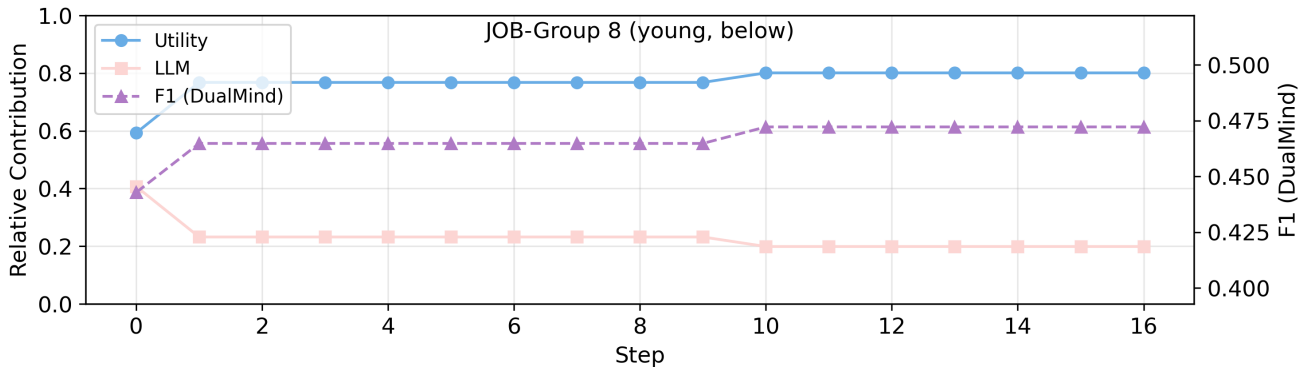
Supplementary Figure 45: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for JOB (old, Low-income) in the JOB dataset.



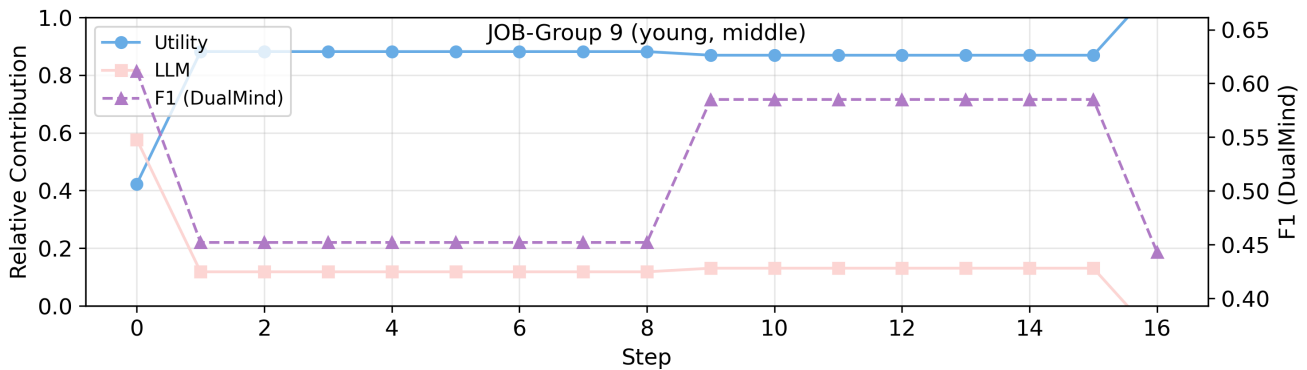
Supplementary Figure 46: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for JOB (old, Middle-income) in the JOB dataset.



Supplementary Figure 47: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for JOB (young, High-income) in the JOB dataset.

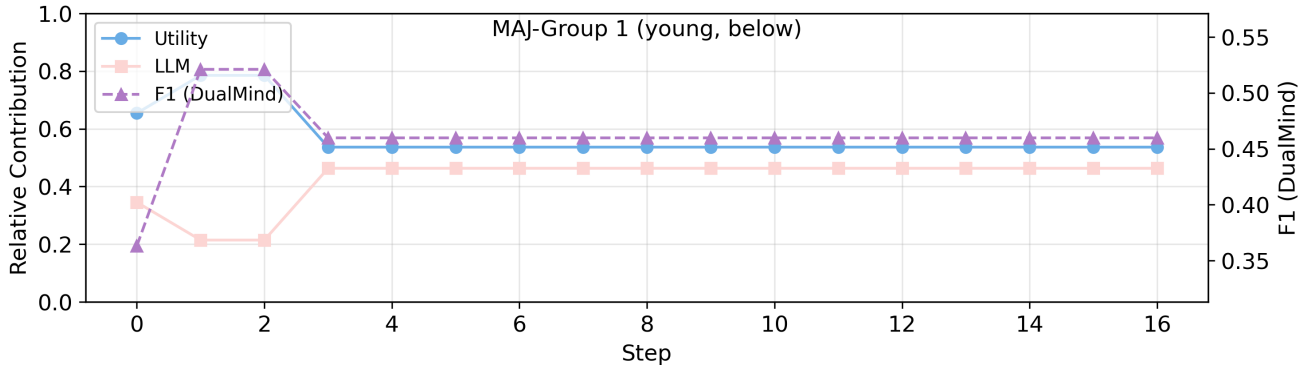


Supplementary Figure 48: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for JOB (young, Low-income) in the JOB dataset.

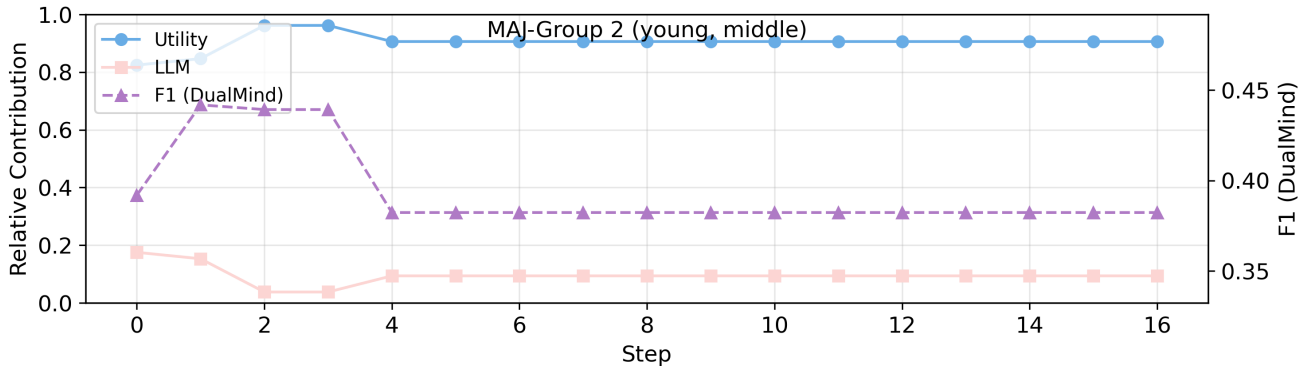


Supplementary Figure 49: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for JOB (young, Middle-income) in the JOB dataset.

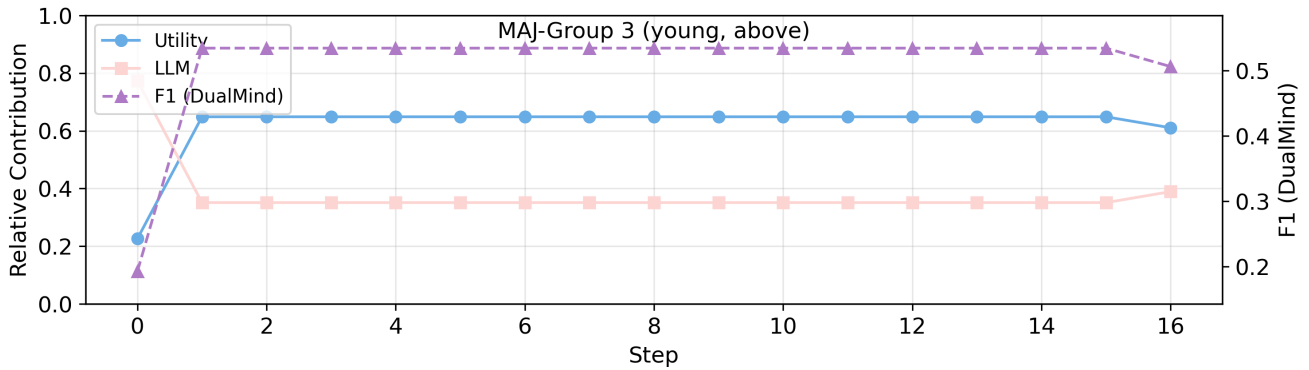
5.6.4 MAJ



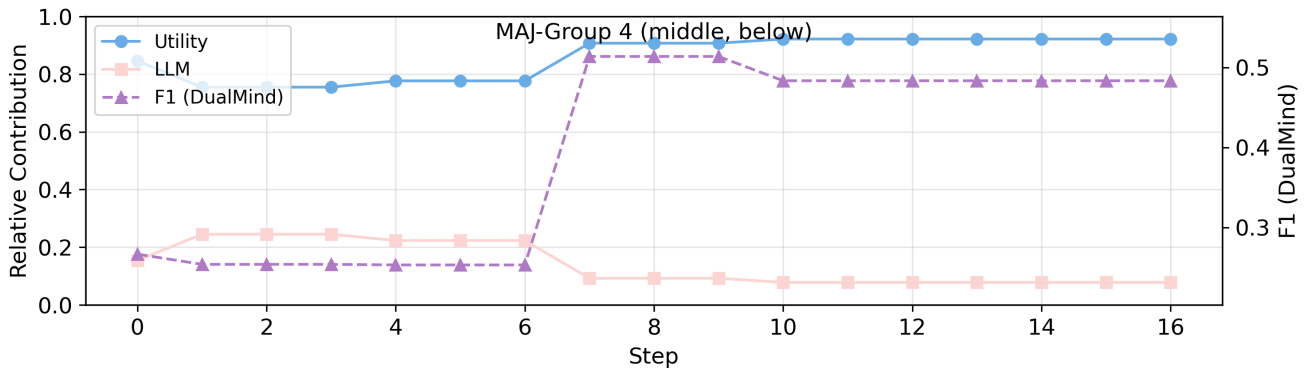
Supplementary Figure 50: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for MAJ (young, Low-income) in the MAJ dataset.



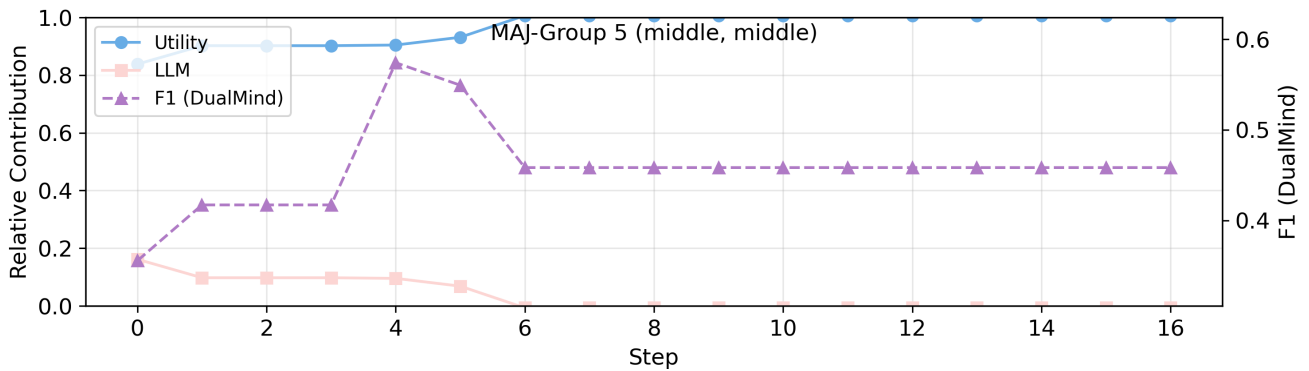
Supplementary Figure 51: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for MAJ (young, Middle-income) in the MAJ dataset.



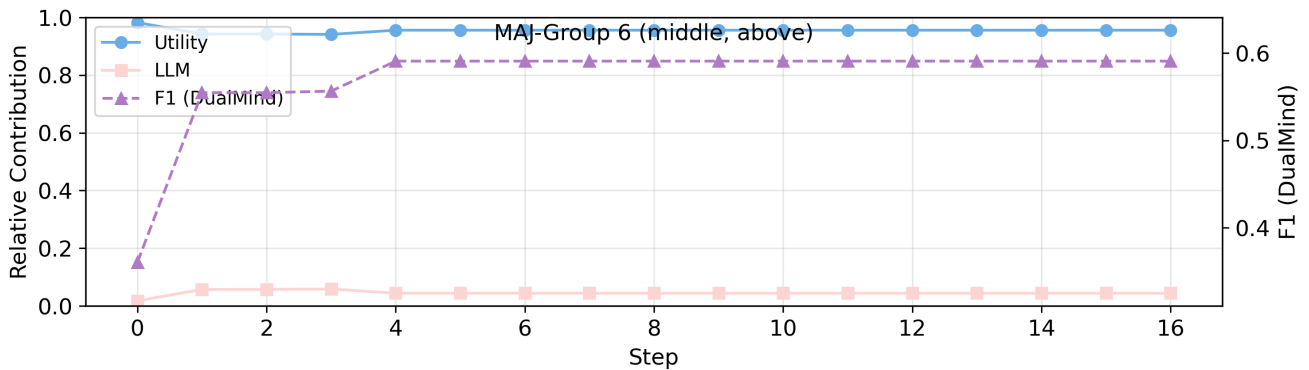
Supplementary Figure 52: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for MAJ (young, High-income) in the MAJ dataset.



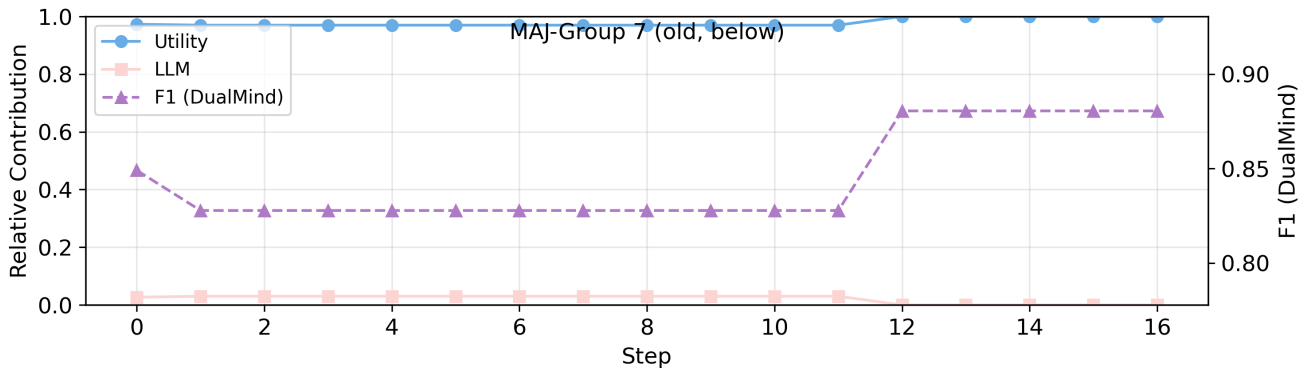
Supplementary Figure 53: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for MAJ (middle, Low-income) in the MAJ dataset.



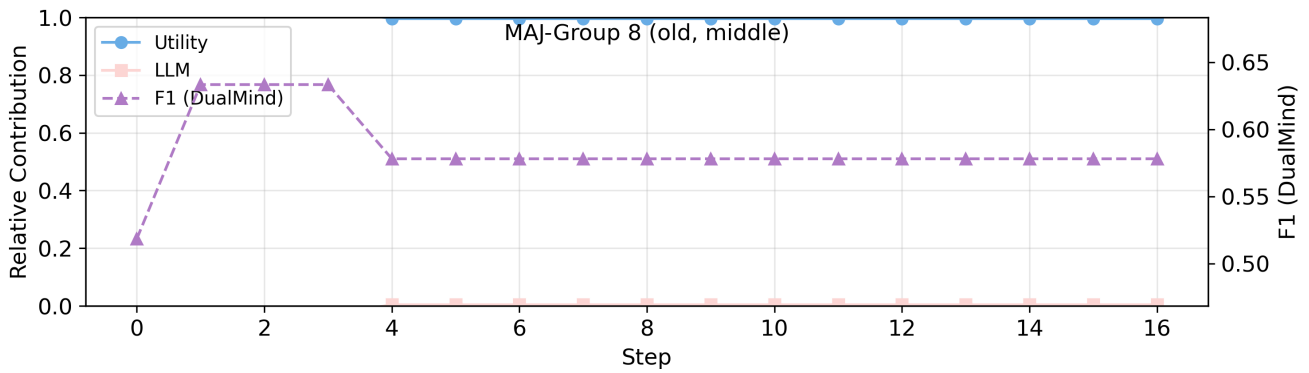
Supplementary Figure 54: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for MAJ (middle, Middle-income) in the MAJ dataset.



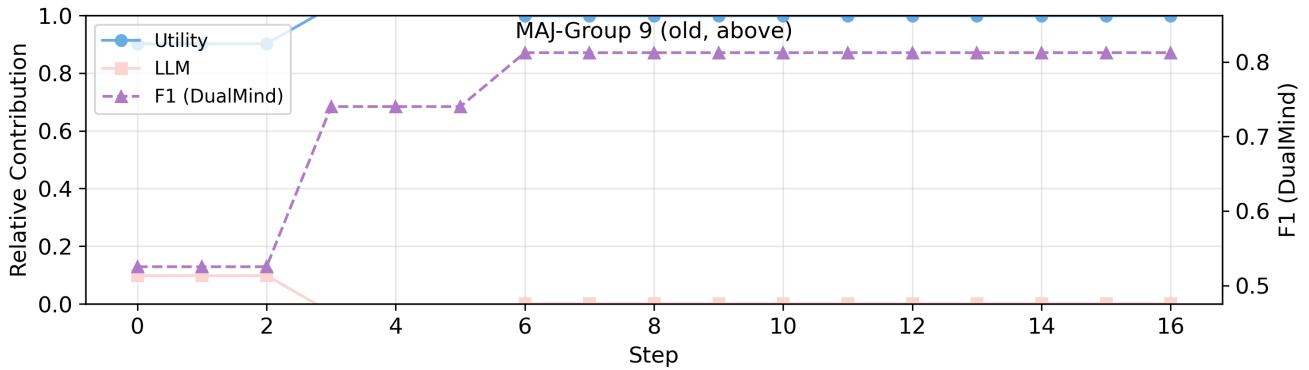
Supplementary Figure 55: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for MAJ (middle, High-income) in the MAJ dataset.



Supplementary Figure 56: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for MAJ (old, Low-income) in the MAJ dataset.

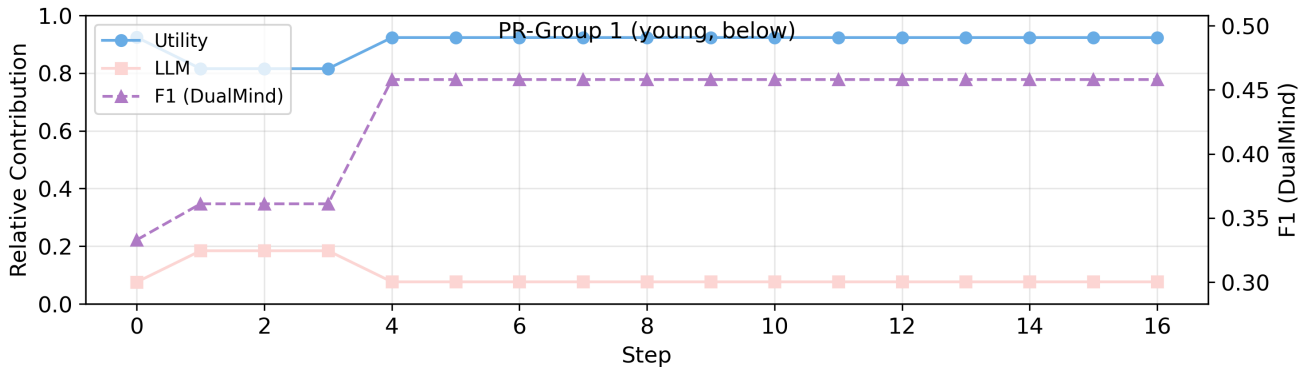


Supplementary Figure 57: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for MAJ (old, Middle-income) in the MAJ dataset.

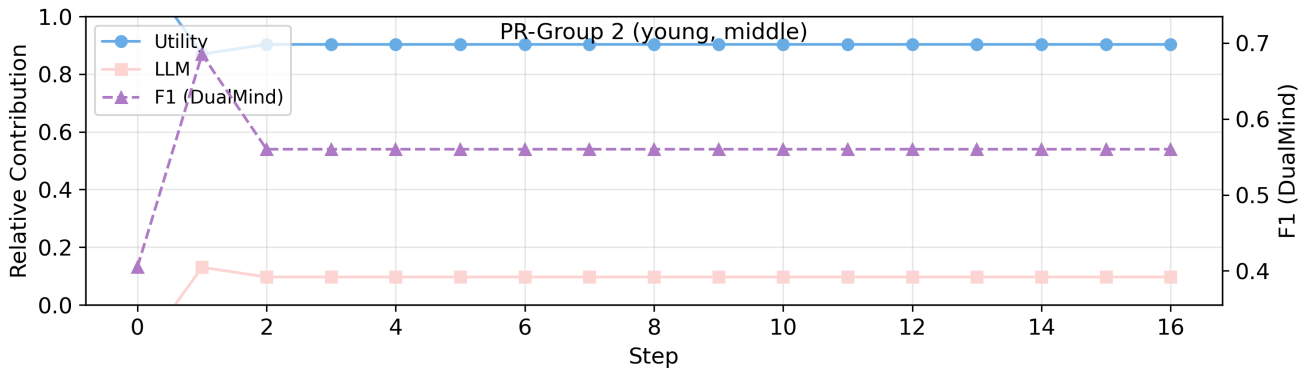


Supplementary Figure 58: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for MAJ (old, High-income) in the MAJ dataset.

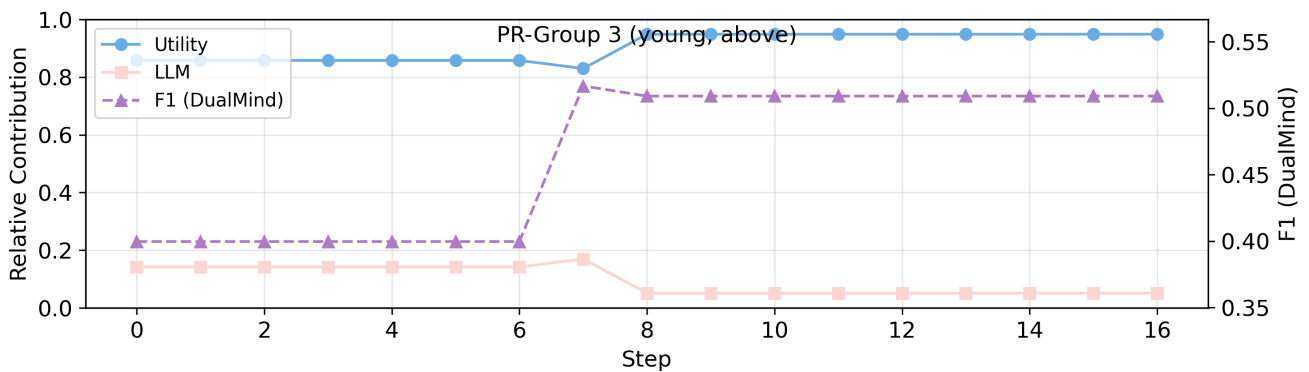
5.6.5 PR



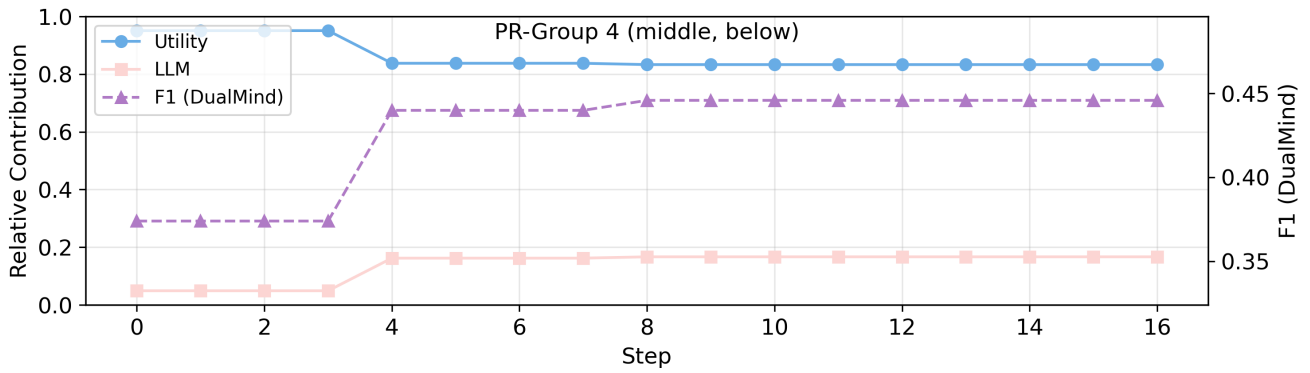
Supplementary Figure 59: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for PR (young, Low-income) in the PR dataset.



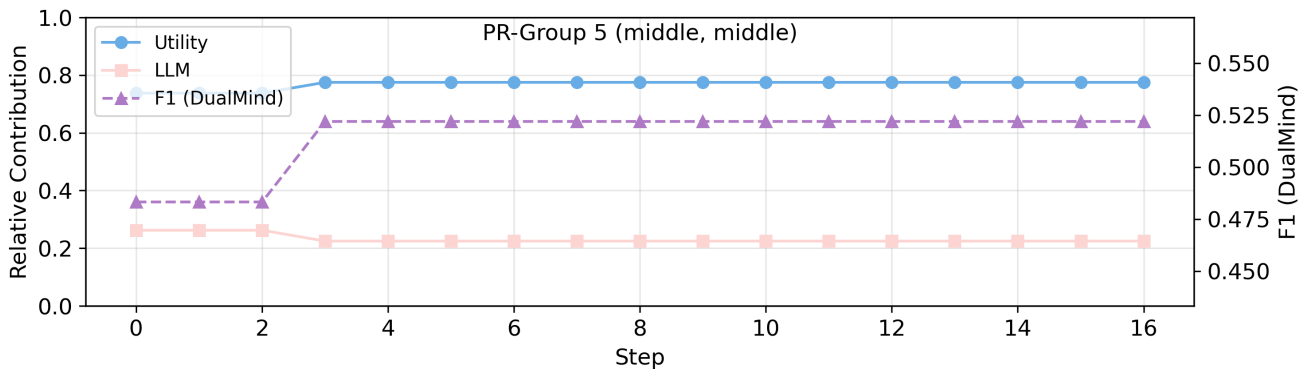
Supplementary Figure 60: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for PR (young, Middle-income) in the PR dataset.



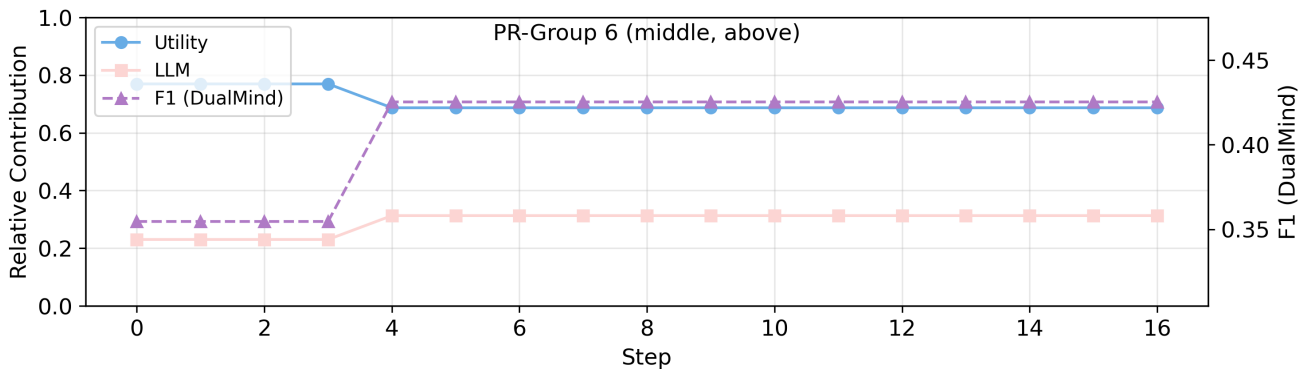
Supplementary Figure 61: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for PR (young, High-income) in the PR dataset.



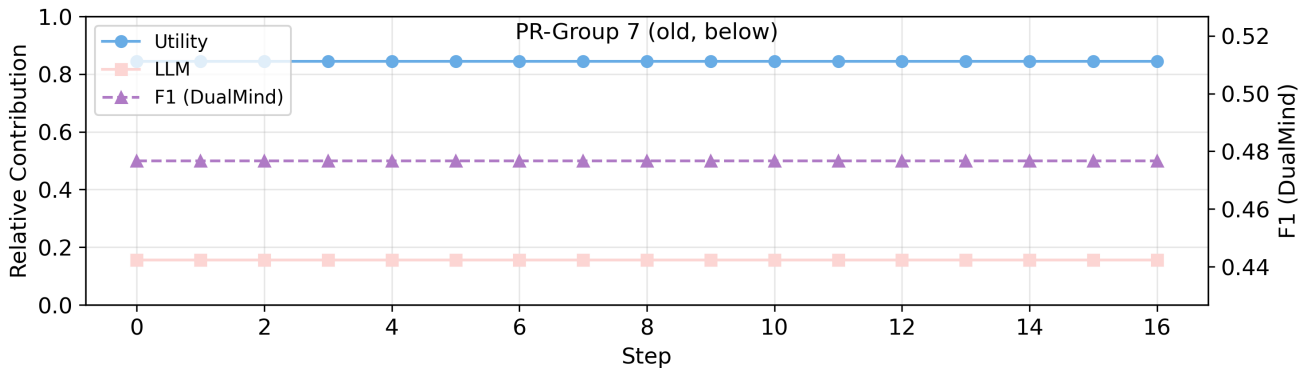
Supplementary Figure 62: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for PR (middle, Low-income) in the PR dataset.



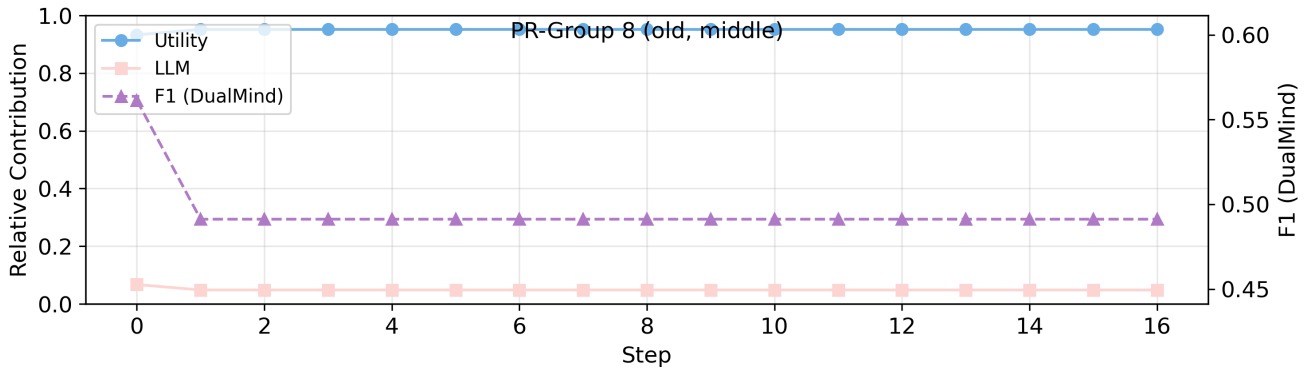
Supplementary Figure 63: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for PR (middle, Middle-income) in the PR dataset.



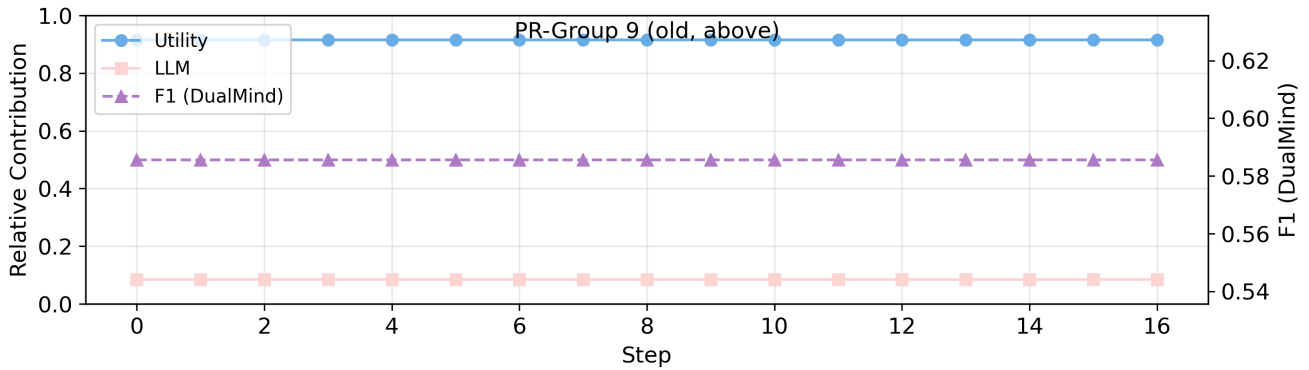
Supplementary Figure 64: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for PR (middle, High-income) in the PR dataset.



Supplementary Figure 65: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for PR (old, Low-income) in the PR dataset.

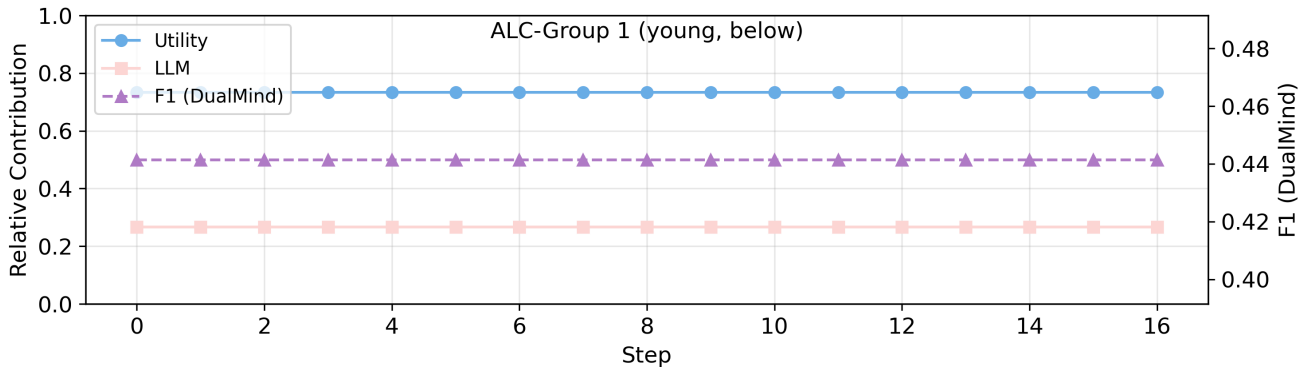


Supplementary Figure 66: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for PR (old, Middle-income) in the PR dataset.

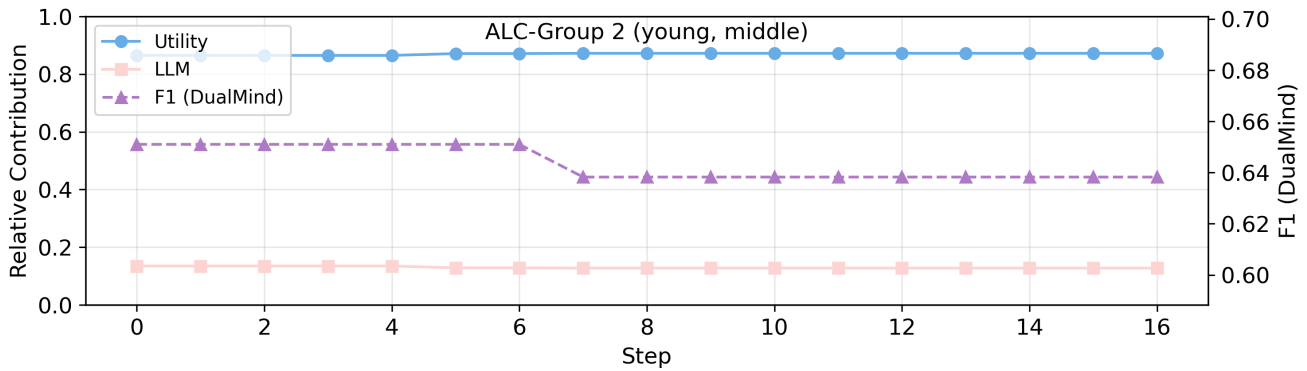


Supplementary Figure 67: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for PR (old, High-income) in the PR dataset.

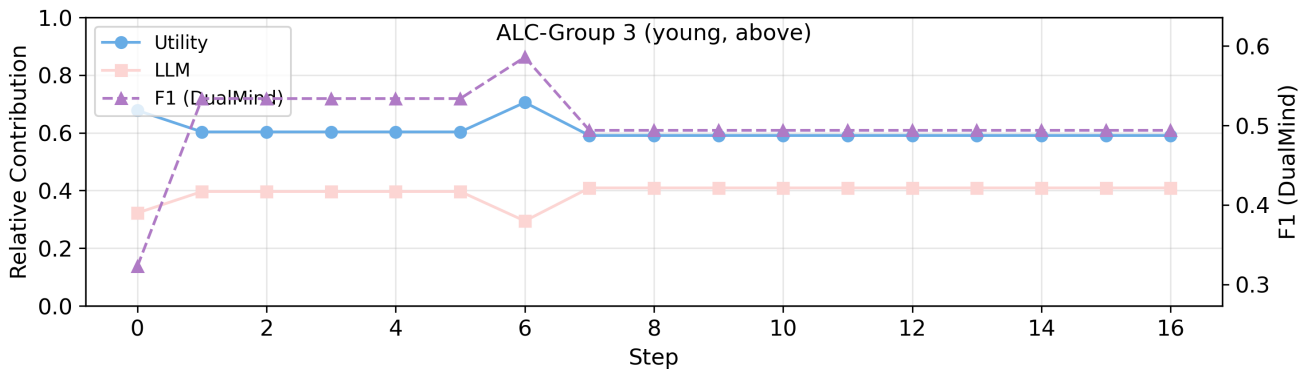
5.6.6 ALC



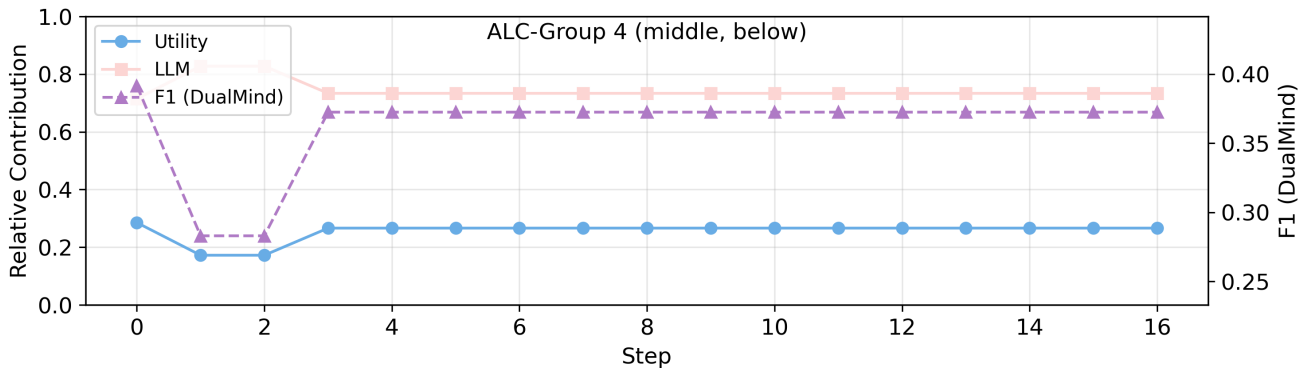
Supplementary Figure 68: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for ALC (young, Low-income) in the ALC dataset.



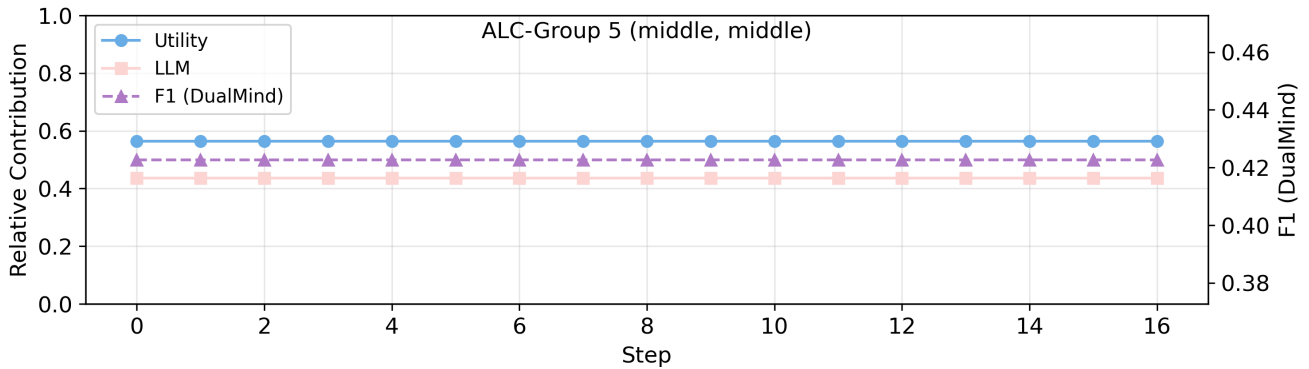
Supplementary Figure 69: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for ALC (young, Middle-income) in the ALC dataset.



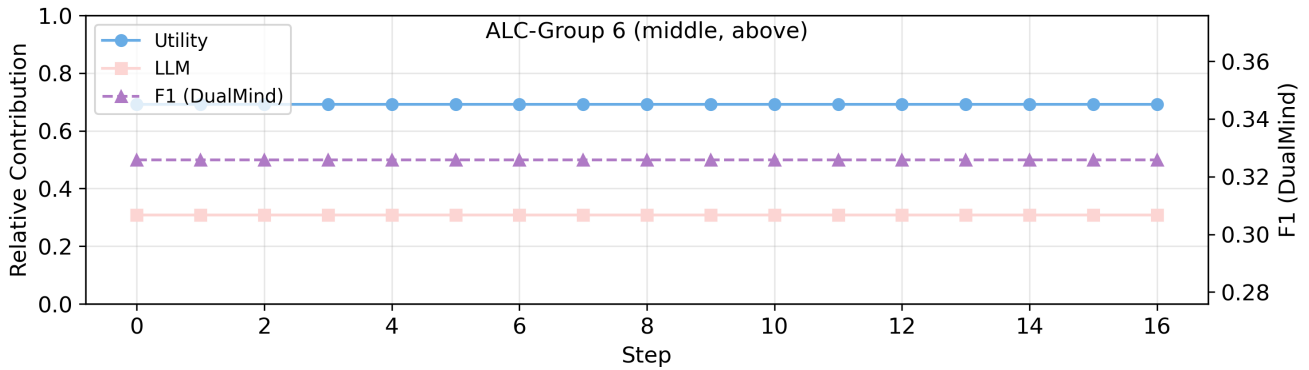
Supplementary Figure 70: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for ALC (young, High-income) in the ALC dataset.



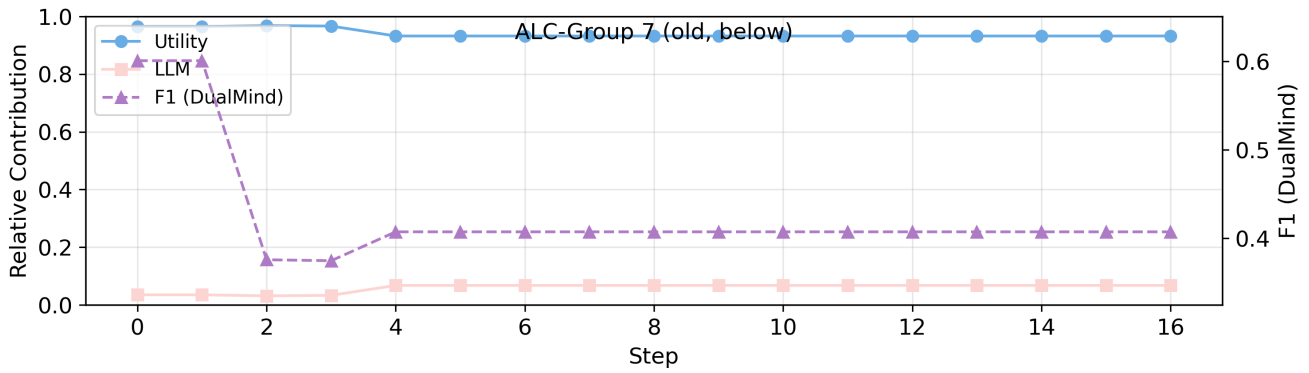
Supplementary Figure 71: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for ALC (middle, Low-income) in the ALC dataset.



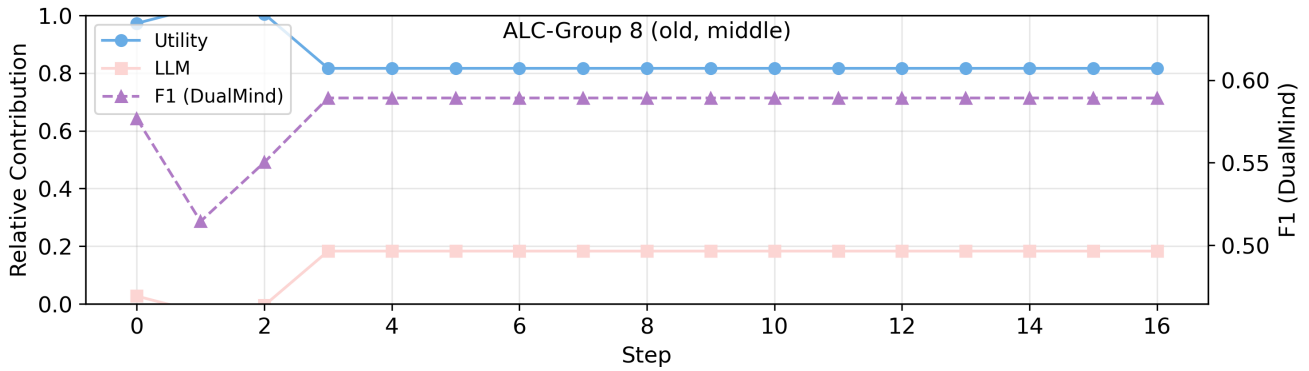
Supplementary Figure 72: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for ALC (middle, Middle-income) in the ALC dataset.



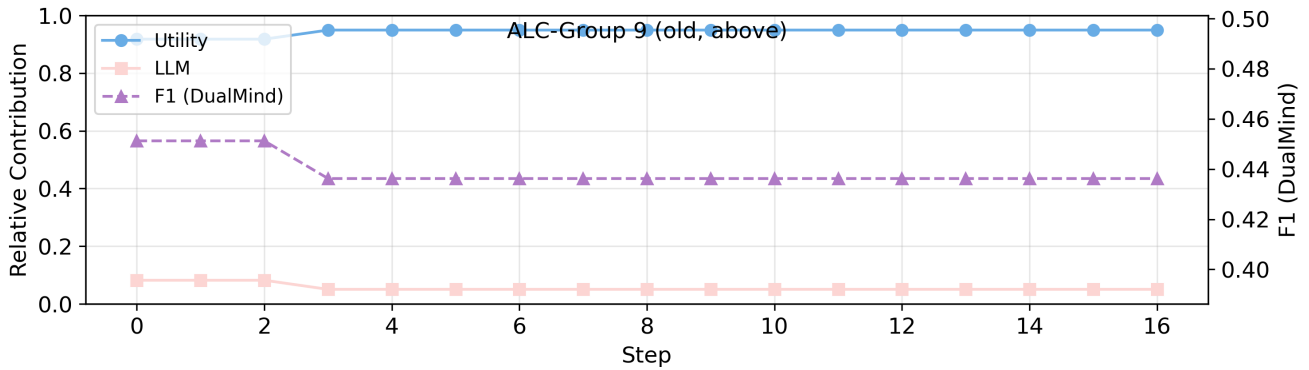
Supplementary Figure 73: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for ALC (middle, High-income) in the ALC dataset.



Supplementary Figure 74: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for ALC (old, Low-income) in the ALC dataset.



Supplementary Figure 75: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for ALC (old, Middle-income) in the ALC dataset.



Supplementary Figure 76: Evolution of the F1 score of the *DualMind* alongside utility-knowledge contributions for ALC (old, High-income) in the ALC dataset.

6 Ablation Study

We conducted an ablation study on the update signals used in *DualMind Gradient*, by varying whether the knowledge set and utility functions were updated using purely textual signals or mixed signals that combine textual feedback with numerical signals. As shown in Supplementary Table 8, the best performance on both the Vaccine and Travel Mode tasks is achieved when the knowledge set is updated with textual signals while the utility functions are updated with mixed signals, yielding the highest F1/Acc on Vaccine (0.662/0.659) and Travel Mode (0.708/0.710).

Supplementary Table 8: Ablation Study of *DualMind Gradient*.

Knowledge Set	Utility Functions	Vaccine		Travel Mode	
		F1	Acc	F1	Acc
Textual	Textual	0.634	0.629	0.688	0.688
Mix	Textual	0.659	0.654	0.614	0.620
Mix	Mix	0.651	0.649	0.683	0.685
Textual	Mix	0.662	0.659	0.708	0.710