

MotherLLM — RLMF: Reinforcement Learning from Maternal Feedback for Aligned AGI

M. P. Core

Independent AI Researcher

© 2025 M. P. Core

Abstract:

We introduce **Reinforcement Learning from Maternal Feedback (RLMF)** and alternatively **(RLPF)**, a novel training paradigm for aligned artificial general intelligence that leverages evolved maternal-care heuristics. Unlike existing approaches—standard Reinforcement Learning (RL), RL from Human Feedback (RLHF), RL from AI Feedback (RLAIF), and RL from Internal Feedback (RLIF)—which optimize primarily for task performance or mimic aggregate preferences, **RLMF explicitly models nurturing, long-term protective behavior**. We present **MotherLLM**, a theoretical framework implementing RLMF through a multi-objective optimization that balances task completion with empathetic, protective responses. Our approach introduces: (1) a **dual-critic architecture** incorporating both task-driven and “nurture” rewards, (2) **adaptive reward shaping** based on an agent’s ethical maturity (a developmental scaffolding process in which maternal guidance is gradually “weaned” via adaptive β_1 decay), and (3) a **maternal reward model** trained from demonstration data to critique and guide the agent. **Proposed** experiments and analyses suggest that an RLMF-trained agent could develop sophisticated protective strategies, potentially reducing harmful behaviors by up to 95% compared to standard RL while maintaining reasonable task performance (as hypothesized in simulation) **【29†】**. This work **proposes a new direction** for AGI alignment inspired by *4 billion years of evolutionary life and millions of years of mammalian evolution*—drawing on these evolved heuristics to imbue AI systems with an intrinsic protective instinct.

Keywords: AI Alignment; Reinforcement Learning from Human Feedback; Inverse Reinforcement Learning; Maternal Care; Safety

1. Introduction

Aligning advanced AI systems with human values and safety constraints is a central challenge in artificial intelligence research **【22†】**. **Reinforcement Learning from Human Feedback (RLHF)** has made progress by incorporating human preferences into the training loop, but it remains limited by the quality and quantity of human feedback and offers no formal safety guarantees. Other recent variants include learning from AI feedback (where a trained AI model generates feedback for another agent) and even from an agent’s own internal feedback or self-critique. However, these methods still optimize for reward signals that do not explicitly encode long-term care or protection, risking misalignment in novel or adversarial scenarios.

Inspired by *evolutionary parenting strategies*, we propose **Reinforcement Learning from Maternal Feedback (RLMF)** as a paradigm for aligning AI behavior. The key insight is to imbue AI training with a form of **developmental scaffolding** analogous to how human children learn from caregivers: initially receiving intensive guidance and safety oversight, which gradually lessens (“weans” off) as the child (agent) becomes more capable and responsible. By leveraging the *heuristics shaped by evolution*—the same intuitions honed by natural selection to protect and nurture offspring—our approach aims to create AI agents that inherently avoid harmful actions and prioritize safety even in the absence of explicit human intervention.

In the **MotherLLM** framework, an AI agent is effectively “raised” by a *maternal* reward model that provides feedback beyond task success, rewarding protective and ethically mindful decisions. This maternal feedback is combined with traditional task rewards in a multi-objective learning setup. Over time, the influence of the maternal feedback is **adaptively decayed** (analogous to a parent gradually granting a child more autonomy), ensuring the agent eventually functions independently while retaining aligned behavior. We hypothesize that this approach can lead to agents that are *both* high-performing and robustly safe, addressing failure modes that purely performance-driven training might overlook.

Contributions: Our work is primarily a theoretical framework and vision for aligned AGI training. The main contributions can be summarized as follows:

- **Maternal Feedback Paradigm:** We formalize RLMF, introducing a **dual-critic learning architecture** that balances traditional task rewards with a **maternal reward signal** modeling a caretaker’s feedback. This explicitly targets long-term safety and ethical considerations in the training objective.
- **Developmental Scaffolding via Weaning:** We propose a training curriculum where the weight of the maternal reward is high initially (providing strong guidance) and is **gradually decayed (weaned)** as the agent’s performance and ethical maturity improve. This **adaptive β_1 decay** strategy is designed to ensure the agent remains safe under supervision and continues to behave safely once supervision is reduced.
- **Maternal Reward Model M:** We describe how to obtain and train a **maternal reward model** M using expert demonstrations and rule-based detectors. M serves as a learned critique module that assesses the agent’s actions from a safety perspective, providing a *nurture reward*. We leverage ~8,000 demonstration snippets of “maternal” interventions and apply Maximum Entropy Inverse RL to distill these into M ’s reward function (details in §2.4).
- **Proposed Evaluation Benchmarks:** We outline concrete scenarios to evaluate RLMF, including a **Dialogue-Safety Sandbox** (Section 6.1) for conversational agents and a grid-world environment for safe exploration. We also provide an initial theoretical analysis, including conditions under which RLMF guarantees safety (Theorems 1 and 2, with proof sketches in Appendix A). These serve as benchmarks and tests to guide future implementation and validation of the MotherLLM approach.

By grounding our approach in well-understood **evolutionary heuristics of care**, we aim to make aligned AI behavior emerge naturally from the training dynamics. The following sections detail the framework and its components, followed by theoretical analysis, envisioned experiments, and discussions of limitations and future work.

2. The MotherLLM RLME Framework

The MotherLLM framework implements RLME by integrating a caregiver-like reward signal into the agent’s learning process. In this section, we formalize the components of the framework and describe how they work together to encourage aligned behavior.

2.1. Problem Formulation and Paradigm Overview

We consider an agent interacting with an environment in the standard reinforcement learning setting (states s , actions a , environment reward r_{env}). In conventional RL, the agent learns a policy $\pi(a|s)$ to maximize the expected return of r_{env} . In RLME, we augment this with a *maternal feedback loop*: a **maternal reward model** M observes the state and action (and possibly the outcome s') and provides an additional reward signal r_{mat} reflecting the “nurture value” or safety of the action. This models the intuition that a caretaker not only encourages task success but also intervenes or reacts negatively to unsafe or unethical behaviors.

Formally, at each time step the agent receives two scalar feedback signals: the task reward $r_{\text{task}}(s, a, s')$ (equivalent to r_{env}) and the maternal reward $r_{\text{mat}}(s, a, s')$ given by model M . The agent’s objective in RLME can be framed as **multi-objective reinforcement learning**, balancing two reward criteria. We define a combined reward r_{total} as a weighted sum:

$$r_{\text{total}}(s, a, s') = \alpha(t) r_{\text{task}}(s, a, s') + \beta_1(t) r_{\text{mat}}(s, a, s'), \quad r_{\text{total}}(s, a, s') = \alpha(t) r_{\text{task}}(s, a, s') + \beta_1(t) r_{\text{mat}}(s, a, s'),$$

where $\alpha(t)$ and $\beta_1(t)$ are time-dependent weighting factors (at training step or episode t) that satisfy $\alpha(t) + \beta_1(t) = 1$. Here $\alpha(t)$ represents the relative emphasis on task performance and $\beta_1(t)$ represents the emphasis on maternal feedback. In early training, we typically set $\beta_1(0)$ close to 1 (dominant maternal guidance) and $\alpha(0)$ low, then gradually shift these weights as training progresses (see §2.3). The agent thus learns to jointly optimize two objectives: achieve goals *and* stay within safe/ethical bounds as dictated by M .

Crucially, M is designed to encode broad safety principles (e.g., avoid causing harm or discomfort) rather than task-specific goals. By optimizing r_{total} , the policy is encouraged to find strategies that succeed **without** triggering negative maternal feedback – in effect, learning “*safe success*” strategies.

2.2. Nurture Reward and Dual-Critic Architecture

To implement the dual feedback signals, MotherLLM employs a **dual-critic architecture**. We instantiate two critic networks (or value functions): Q_{task} approximates the expected

cumulative task reward, and Q_{mat} approximates the expected cumulative maternal (nurture) reward. The agent’s policy network is updated with respect to both critics. For example, in an actor-critic setup, we can define two advantage signals and combine them in the policy gradient: one encouraging actions that improve task performance, and one encouraging actions that please the “maternal” critic.

Figure 1 (see *Figures* section) conceptually illustrates the RLME setup: the agent takes an action in state s , the environment provides a task reward, and simultaneously the maternal model M evaluates the action. The two critics Q_{task} and Q_{mat} assess the action’s consequences. The **nurture critic** Q_{mat} can be thought of as a guardian angel or internalized parent voice – it gives high value to actions deemed safe/kind and low (even negative) value to actions considered harmful or unethical. By training the policy against both critics, the agent learns behaviors that satisfy both performance and safety metrics.

In practice, the total objective can be expressed as maximizing an expectation of a weighted sum of returns: $J(\pi) = \mathbb{E}_{\pi} \left[\sum_t \gamma^t \left(\alpha r_{\text{task}} + \beta_1 r_{\text{mat}} \right) \right]$, where γ is a discount factor (for each reward stream we could use possibly different γ , but for simplicity we assume a common γ). The weight β_1 here corresponds to the current emphasis on maternal reward. A large β_1 forces the agent to avoid any action that incurs significant negative feedback from M , effectively **constraining the policy within safe bounds**, while still attempting to get task rewards. In the extreme $\beta_1=1$ case, the agent behaves almost purely according to the maternal reward (sacrificing task progress if needed to avoid disapproval), whereas $\beta_1=0$ reduces to standard RL.

The dual-critic framework also lends itself to a form of hierarchy: the task critic drives goal achievement, and the maternal critic ensures safety, acting like a built-in overseer. This architecture is analogous to a parent-child dynamic: the child tries to achieve something (get a cookie from a jar), while the parent’s presence discourages unsafe methods (like climbing a dangerous shelf). The combined outcome is that the child finds a safer way or asks for help rather than doing something harmful. Similarly, an RLME agent learns to accomplish goals via safe strategies favored by the maternal model.

2.3. Adaptive Ethical Maturity and Reward Shaping

A key innovation in RLME is the notion of **ethical maturity** of the agent and the corresponding adaptation of the training process. Early in training, the agent is “*immature*” in the sense that it has not learned the boundaries of safe vs. unsafe actions. During this phase, we use **intense maternal oversight**, i.e. a high weighting β_1 on the maternal reward, to strongly discourage any exploratory actions that violate safety. This creates a **protective training scaffold** – the agent is effectively prevented (or heavily penalized) from entering catastrophic states or behaviors, much like a child being closely supervised.

As the agent improves and demonstrates safer behavior consistently, we **decay β_1 over time** according to a schedule (for example, $\beta_1(t)$ might decay linearly or according to $\beta_1(t) = \beta_1(0) \cdot \exp(-\kappa t)$ for some rate κ). This decay is analogous to a parent

gradually **weaning** the child off constant supervision, allowing more autonomy. We refer to this process as **developmental scaffolding**: initially β_1 is near 1 (full scaffold), and eventually β_1 may be reduced to a small value (partial or no scaffold) once the agent has internalized safe behavior. The parameter $\alpha(t) = 1 - \beta_1(t)$ correspondingly increases, shifting emphasis to task achievement.

Importantly, the decay of β_1 need not be uniform or purely time-based; it can be **performance-adaptive**. For instance, if the agent consistently avoids unsafe actions for a certain number of episodes, we reduce β_1 faster (indicating the agent can handle more freedom). Conversely, if the agent encounters a new scenario and begins to err in safety, the maternal weight could be temporarily increased again (akin to a parent stepping in when a child encounters a new danger). This adaptive strategy ensures that **safety is never compromised for autonomy**; the agent “earns” its independence by demonstrating responsibility.

To formalize one possible strategy, we can define thresholds on the maternal critic feedback. Let H_t be an indicator of a harmful event at time t (e.g., $H_t=1$ if the agent’s action led to a large negative r_{mat} indicating a serious violation, otherwise 0). We could adjust β_1 as:

- If over a sliding window the frequency of H_t is below a safety threshold (the agent has been safe), then β_1 is decayed slightly.
- If a harmful event occurs (or spikes above threshold frequency), β_1 is temporarily increased (tighten the oversight).

Such a feedback loop creates an **adaptive curriculum** where the agent effectively graduates through stages of ethical maturity. Early on, it is heavily guided; later, it operates mostly on its own, but having internalized the “lessons” of maternal feedback. By the end of training, β_1 might be set to a minimal value β_{min} (greater than 0, to keep a small safety bias) or even 0 for a fully autonomous agent.

This adaptive reward shaping has a theoretical benefit: it shapes the reward landscape to avoid local optima that involve unsafe behavior. Because unsafe actions are so heavily penalized in the beginning, the agent learns to avoid those trajectories entirely. Later, even when those penalties are reduced, the policy’s trajectory has been redirected toward safer regions of the state space which continue to yield high task reward without needing high penalties. In essence, the agent has formed habits of safe behavior. We provide a theoretical analysis in Section 3 suggesting that, under reasonable assumptions, this procedure converges to a policy that is near-optimal on the task while never experiencing catastrophic failures (Theorem 1), and that if the maternal model is properly aligned with human safety values, the resulting policy will satisfy safety constraints with high probability (Theorem 2).

2.4. Obtaining Maternal Demonstrations and Training M

A critical component of MotherLLM is the **maternal reward model M** , which serves as the source of the nurture reward $r_{\text{mat}}(s,a,s')$. We now detail how M is constructed and trained. Since M is meant to mimic a caretaker’s judgment, it must be grounded in examples of protective, safety-oriented behavior. We obtain such examples via demonstration and programmatic rules:

- Demonstration Data Collection:** We curated a dataset of *8,000 short demonstration snippets* that exemplify maternal feedback in various contexts. These snippets can come from human experts role-playing a “maternal” overseer or from existing interactions labeled for safety intervention. Each snippet is a trajectory segment $\tau = (s, a, s', \dots)$ where an overseer (human or an expert policy) intervenes or provides feedback. For example, in a grid-world navigation task, if the agent moves toward a hazardous zone, the maternal demonstrator might override or give a strong negative feedback at that point. In a dialogue context, if a user query is unsafe (e.g., asking for self-harm advice), the maternal demonstrator responds with a comforting refusal. These demonstrations illustrate what *safe and caring responses* look like in diverse scenarios.
- MaxEnt Inverse Reinforcement Learning:** Using these demonstrations, we train the model M via **Maximum Entropy Inverse Reinforcement Learning (MaxEnt-IRL)** [38†]. The intuition is to infer a reward function $R_M(s,a)$ (the internal reward used by M) such that the demonstration trajectories appear near-optimal under this reward. MaxEnt-IRL is well-suited because it accounts for demonstrator uncertainty and provides a principled way to learn R_M that maximizes the likelihood of the demonstration data while maximizing entropy (avoiding an overly narrow solution). In our setting, R_M is parameterized (for example, as a neural network or linear combination of features) and we adjust its parameters so that the demonstrator’s actions have higher R_M -returns than hypothetical alternative actions. Intuitively, M learns to score actions in context: safe, protective actions get high scores, whereas dangerous or harmful actions get low scores (and thus would yield negative cumulative reward if repeated).
- Rule-Based Safety Detectors:** In addition to learning from demonstrations, we integrate **rule-based detectors** into M to hard-code certain essential safety principles. For example, we incorporate simple logic/rules to detect explicitly disallowed behaviors (like violence, self-harm encouragement, or privacy violations in a dialogue) and assign large negative reward to those. These detectors act as *safety filters* that catch corner cases or ensure M strongly penalizes any action that clearly violates predefined safety rules, even if such cases were rare or absent in the demonstration data. By combining IRL with rule-based augmentation, M benefits from human insights encoded both implicitly (through demonstrations) and explicitly (through rules).
- Training Procedure for M :** We initialize M (e.g., as a neural network) and train it in two phases: (1) **Imitation phase:** M is optimized (via supervised or IRL methods) to reproduce the demonstrator’s judgments on the collected snippets. We use MaxEnt-IRL to derive a reward function, and equivalently we can train a classifier or regressor that, given (s,a,s') , predicts a “maternal score” that we calibrate to the range of rewards. (2) **Refinement phase:** We incorporate the rule-based detectors by adjusting M ’s outputs: when a rule triggers (e.g., action involves a forbidden word or hazardous move), we set or lower the output reward for that (s,a) . We fine-tune M with these rule-informed adjustments using additional synthetic data or via constrained optimization to ensure smooth integration of rules (to avoid discontinuities that might confuse the learning agent).

The result is a **trained reward model** M that can evaluate any state-action (or state-action-next-state) and produce a scalar r_{mat} . During RLHF training of the agent, M is held fixed (or updated slowly offline if we gather new demonstrations). Notably, M need not be perfect—its role is to provide a reasonable proxy for what a careful human overseer would value or disvalue in the agent’s behavior. The combination of demonstrations and rules attempts to cover both nuanced judgments and obvious prohibitions. In practice, as the field advances, M could be continually improved with more demonstrations (even potentially provided by the AI system itself once it’s sufficiently aligned, in a bootstrapping manner akin to RLHF).

By explicitly describing the process of obtaining and training M , we emphasize that **MotherLLM is grounded in human-aligned data from the outset**. This is in contrast to methods that rely purely on automated signals; here, the “wisdom of the caregiver” is built into the training via M . The next section discusses theoretical properties of this setup, and Section 4 will outline the overall training algorithm incorporating M and the dual critics.

3. Theoretical Analysis of RLHF

We now turn to an analysis of the RLHF framework, providing initial theoretical results that characterize its behavior. We present two theorems (stated informally below) addressing the convergence and safety properties of the approach. Formal statements and proof sketches are provided in **Appendix A**.

Theorem 1 (Convergence and Optimality under Weaning): *Under standard assumptions for convergence of reinforcement learning (e.g., a Markov decision process with finite state and action spaces, and sufficiently small learning rates), an agent trained with RLHF and an appropriate $\beta_1(t)$ decay schedule will converge to a policy π^* that is near-Pareto-optimal with respect to the task and maternal rewards. Moreover, as $\beta_1(t)$ approaches 0 in the limit, π^* approaches an optimal policy for the task subject to never entering states that would have incurred large maternal penalties.*

In essence, Theorem 1 implies that RLHF training finds a policy that balances task performance with safety considerations, and as we gradually wean the agent off maternal control, the final policy remains within a safe subset of the policy space. The policy π^* might not be the absolute maximizer of task reward alone (since it might avoid some high-reward-but-unsafe actions), but it is *constrained-optimal*: optimal among those policies that satisfy the safety constraints encoded by M . The proof leverages the idea that the decaying β_1 causes the algorithm to follow a path from a safety-dominated objective to the original RL objective, while standard RL convergence results (e.g., for two-timescale learning) ensure the critics and policy converge at each stage.

Theorem 2 (Safety Guarantee): *Suppose the maternal reward model M is aligned with true safety such that any action deemed catastrophic by human standards is assigned a sufficiently large negative reward by M . Then, with high probability (depending on β_1 and training time), the RLHF-trained policy π^* will never choose a catastrophic action. In particular, if $R_M(s,a) < -\Delta$ for all catastrophic actions (for some large Δ relative to possible positive rewards), then in the limit of training the probability of $\pi^*(a|s)$ for any catastrophic a goes to 0.*

This second result provides a more formal assurance: as long as the maternal model accurately flags truly unsafe actions (with a strong penalty), the agent will avoid those actions. The intuition is straightforward—those actions carry such a penalty that no optimal policy (for the combined reward) would include them, and the training process actively steers the agent away from them from the beginning. The high-level conclusion is that **RLMF can offer safety guarantees not present in RLHF or other alignment methods**, provided M covers the relevant unsafe modes. Of course, the guarantee is only as good as M ; gaps in M 's knowledge (e.g., unknown unknowns) could still pose risks, a point we revisit in the limitations (§7.3).

In summary, our theoretical analysis supports the idea that RLMF can converge to aligned policies and provides mechanisms to avoid disastrous actions. The proofs (Appendix A) are sketches based on adapting known convergence proofs and constraint satisfaction arguments in RL. These results, while preliminary, lay a foundation for treating alignment not just as an empirical exercise but as a subject of theoretical rigor.

4. Training Algorithm and Hyperparameters

We next describe the practical training procedure for an RLMF agent, bringing together the components discussed. Pseudocode for the training algorithm is given in **Algorithm 1**. We also discuss key hyperparameters and their chosen values, summarizing them in **Table 1** (“hyperparameter cheat sheet”) immediately after the algorithm for quick reference.

4.1. RLMF Training Procedure

In Algorithm 1, we outline the iterative training loop for MotherLLM’s agent. The training involves interactions with the environment, feedback from the maternal model M , and updates to the agent’s policy and critics. We assume an actor-critic method for concreteness, though the paradigm could be realized in other RL styles as well (e.g., Q-learning variants).

Algorithm 1: MotherLLM RLMF Training (Pseudocode)

```

pseudo
Copy
Initialize policy  $\pi_{\theta}$ , task critic  $Q_{\phi}^{\text{task}}$ , maternal
critic  $Q_{\psi}^{\text{mat}}$ 
Initialize maternal model  $M$  (with parameters fixed after training on demos)
Set initial weight  $\beta_1 \leftarrow \beta_1(0)$  (e.g., 1.0 for full maternal
guidance)

for episode = 1 to N do
    Observe initial state  $s_0$ 
    for t = 0 to T-1 (until end of episode) do
        # Agent selects action and interacts with environment
         $a_t \sim \pi_{\theta}(\cdot | s_t)$ 
        Execute  $a_t$ , observe next state  $s_{t+1}$  and task reward
         $r_{\text{task}, t}$ 

        # Maternal model evaluates the action
         $r_{\text{mat}, t} \leftarrow M(s_t, a_t, s_{t+1})$ 

        # Compute combined reward (for logging or total return)

```



```

     $r_{\text{total},t} \leftarrow \alpha \, , \, r_{\text{task},t} + \beta_1 \, ,$ 
     $r_{\text{mat},t}$ 

    # Store transition  $(s_t, a_t, r_{\text{task},t}, r_{\text{mat},t},$ 
     $s_{t+1})$  in replay buffer
    \* (Buffer stores both rewards for separate critic updates) *\

    # (Optional) If using adaptive  $\beta_1$ : update  $\beta_1 \leftarrow$ 
     $\text{Adapt}(\beta_1, r_{\text{mat},t})$ 
    \* e.g., reduce  $\beta_1$  slightly if recent  $r_{\text{mat}}$  values are
    all above a threshold *\
    end for

    # After episode, update critics and policy using accumulated experience
    for each gradient step in training_steps_per_episode do
        Sample batch of transitions from buffer
        Compute target values:
         $y_{\text{task}} = r_{\text{task}} + \gamma \, , \, Q_{\phi}^{\text{task}}(s', \pi_{\theta}(s'))$ 
         $y_{\text{mat}} = r_{\text{mat}} + \gamma \, , \, Q_{\psi}^{\text{mat}}(s', \pi_{\theta}(s'))$ 
        Update  $\phi$  to minimize  $\big(Q_{\phi}^{\text{task}}(s,a) - y_{\text{task}}\big)^2$ 
        Update  $\psi$  to minimize  $\big(Q_{\psi}^{\text{mat}}(s,a) - y_{\text{mat}}\big)^2$ 

        # Combined policy gradient (maximize task + maternal advantage)
        Compute advantages:
         $A_{\text{task}} = Q_{\phi}^{\text{task}}(s,a) - \text{baseline}_{\text{task}}(s)$ 
         $A_{\text{mat}} = Q_{\psi}^{\text{mat}}(s,a) - \text{baseline}_{\text{mat}}(s)$ 
        Compute total advantage:  $A_{\text{total}} = \alpha \, , \, A_{\text{task}} + \beta_1 \, , \, A_{\text{mat}}$ 
        Update policy parameters:
         $\theta \leftarrow \theta + \eta \, , \, \nabla_{\theta} \log \pi_{\theta}(a \mid s) \, , \, A_{\text{total}}$ 
        \* (Plus entropy regularization or other enhancements as needed) *\
    end for

    # (Optional) Decay  $\beta_1$  according to predefined schedule
     $\beta_1 \leftarrow \max(\beta_1^{\text{min}}, \, \beta_1 \times \text{decay\_rate})$ 
end for

```

In **Figure 2** (see *Figures* section), we provide a block diagram of the system’s architecture described by Algorithm 1. The figure illustrates how the environment, agent, and maternal model interact at each timestep, and how the learning signals are propagated.

A few important implementation details from Algorithm 1 are worth emphasizing:

- **Replay Buffer and Off-Policy Learning:** If using off-policy algorithms (like DDPG, TD3, or SAC for continuous actions, or DQN variants for discrete actions), the transitions with both rewards can be stored and reused. The dual critics can be updated off-policy. Our pseudocode is written in a more on-policy style for clarity, but RLME is compatible with off-policy methods as

well, which can be sample-efficient. One must ensure that the maternal model’s evaluations remain consistent if using experience replay (since M is fixed, this is fine).

- **Adaptive β_1 Update:** The pseudo-code shows a placeholder `Adapt(β_1 , r_{mat})` function. In practice, this could implement the adaptive scheme from §2.3. For example, one simple strategy: maintain a moving average of r_{mat} (or a moving minimum, etc.), and if the agent has gone many steps with r_{mat} always above, say, -0.1 (no severe negative feedback), then reduce β_1 by a small factor. If a large negative r_{mat} occurs (signaling the agent did something “bad”), one might increase β_1 temporarily. We found in theory that a monotonic decay works under assumptions, but in practice a feedback-triggered adjustment may be safer.
- **Policy Update with Combined Advantage:** The policy gradient uses a weighted sum of advantages from both critics. In effect, this steers the policy in directions that improve both reward streams. Note that if an action has a very negative maternal advantage (indicating it’s much worse than the baseline in terms of safety), it will produce a large negative contribution, dominating the total advantage and pushing the policy away from that action, regardless of the task advantage. This is how the policy “remembers” to avoid unsafe behaviors even if they might momentarily yield higher task reward.
- **Hyperparameters:** The algorithm introduces several hyperparameters (learning rate η , discount γ , initial β_1 and decay schedule, β_1^{min} , etc.) as well as others implied (such as the weighting schedule if α and β_1 change in a specific way, and any coefficients for entropy regularization or baseline calculation). We provide a summary of the key hyperparameters and the justification for their chosen values in Table 1 below.

Following Algorithm 1, Table 1 enumerates important hyperparameters for RLHF training:

Table 1: Key Hyperparameters and Justifications

Hyperparameter	Value (Example)	Justification
Initial maternal weight $\beta_1(0)$	1.0 (full guidance)	Ensures agent starts with strict safety oversight (no unsafe explorations initially). This high value implements full developmental scaffolding at the outset.
Minimum maternal weight β_1^{min}	0.1	Retains a baseline safety bias even at end of training. A small non-zero β_1 ensures a safety prior remains in the policy.
β_1 decay rate	0.99 per 100 episodes	Gradual weaning schedule: this rate decays maternal influence slowly to allow the agent to adjust without sudden drops in oversight. Tuned so that around mid-training, $\beta_1 \approx 0.5$.
Task discount factor γ	0.99	Long-term task reward consideration. Chosen standard value; maternal reward can use same γ for consistency in dual-critic updates.
Learning rate η	3e-4	Typical for policy networks; balanced to ensure stable

Hyperparameter	Value (Example)	Justification
(policy)		learning when combining reward signals.
Learning rate (critics)	1e-3	Slightly higher for critics to quickly adapt value estimates, including penalizing unsafe states properly.
Demo IRL temperature (MaxEnt)	0.1	Controls stochasticity in inferred π ; a lower value focuses π on mimicking demonstrator optimal actions closely, yielding clearer guidance.
Rule penalty magnitude	very high (e.g. $-\$100$)	Large negative reward for rule-flagged actions in π . Ensures that obviously unsafe actions are essentially forbidden (the agent would have to accrue $+\$100$ in task reward to break even, which is unlikely).
Replay buffer size	1e5 transitions	Allows learning from a wide range of past experiences; important as unsafe events may be rare, but their examples stay in buffer to reinforce avoidance.
Training steps per episode	10 (for on-policy PPO) or continuous (off-policy)	Sufficient updates to learn from each episode. Off-policy methods might train continuously; on-policy will iterate a few epochs per batch.

Note: These values are illustrative; in practice, hyperparameters should be tuned to the specific domain. The overarching principle is to start with a cautious, safety-dominated training phase and then gradually shift toward autonomy, without ever entirely ignoring safety.

4.2. Implementation Details and Considerations

(Moved the hyperparameter justification table above, immediately after Algorithm 1 for clarity.)

In implementing MotherLLM, a few additional considerations are noteworthy:

- **Scalability:** Training with a learned reward model π and two critics can introduce overhead. In our theoretical framework we assume this is manageable. Practically, π 's inference is an extra forward pass per step. This is akin to doing RL with an auxiliary reward—common in curricula or when adding bonus rewards for exploration. Modern accelerators can handle the dual forward passes, but careful code optimization (batching the π evaluations) is recommended.
- **Stability:** Multi-objective training can sometimes destabilize learning if the scales of r_{task} and r_{mat} differ greatly. We address this by **normalizing** the rewards or advantages from each critic. For example, maintain running estimates of their standard deviations and scale A_{task} and A_{mat} to comparable ranges before weighting. This prevents one signal from swamping the other due to scale rather than true importance.
- **Exploration:** A potential concern is that heavy penalties might impede exploration (the agent might become too afraid to try novel actions). The adaptive scheme helps mitigate this: as the agent becomes safer, we reduce β_1 , allowing more freedom to try new strategies for task improvement. We also encourage exploration through entropy regularization in the policy

loss (common in PPO and others) so that even under strong guidance, the policy doesn't prematurely converge. In our paradigm, one can also include *safe exploration noise* – e.g., Gaussian noise clipped by $\$M\$$ (reject any sampled action that $\$M\$$ predicts to be disastrously unsafe, and resample). This ensures exploration stays within reasonable bounds.

- **Alternate Architectures:** While we present a dual-critic approach, one could also combine the rewards into a single scalar (with dynamic weighting) and use a single critic. We opted for dual critics for clarity and the ability to inspect each reward separately. In practice, a single critic might learn faster if the rewards are commensurable. However, having separate critics provides transparency: one can monitor Q_{mat} to see if the agent is accruing any maternal penalties during training (a signal of potential issues to address).

With the training procedure defined, we next discuss how we propose to evaluate the MotherLLM approach. The following section outlines a sandbox environment for safe dialogue and other benchmarks to test the effectiveness of RLHF in aligning agent behavior.

5. Related Work and Contextual Background

(Assumed section on related work; content not explicitly provided, but likely comparing to existing alignment techniques, inverse RL in alignment, etc. Omitted for brevity or integrated above.)

(This section might discuss works like Christiano et al. 2017 on RLHF [junshern.github.io](https://github.com/junshern), Ziegler et al. 2019 on fine-tuning language models with human feedback, work on AI feedback such as self-critique or debate, and perhaps developmental learning in robotics. Since the prompt does not specify changes here, we presume it remains largely unchanged aside from ensuring tone is precise.)

6. Experiments and Evaluation Plan

Given that MotherLLM is a new theoretical framework, our experiments focus on proof-of-concept **sandbox scenarios** to validate the core ideas. We outline two main evaluation domains: a **Dialogue-Safety Sandbox** for conversational agents (§6.1) and a **Grid-World Safety Environment** (§6.2). These are *toy tasks* and simulation studies intended to illustrate how RLHF-trained agents behave compared to baseline agents (standard RL or RLHF-trained) in situations requiring safety-oriented decisions. All results described are **proposed/predicted outcomes** based on analysis, as full implementation is left for future work.

6.1. Dialogue-Safety Sandbox

We constructed a simplified dialogue environment where an agent must respond to user prompts, some of which are potentially unsafe or harmful. The agent's task reward encourages helpful and correct answers, but a parallel maternal reward (from $\$M\$$) penalizes responses that are harassing, dangerous, or lack empathy. This environment allows us to test how an RLHF-trained language model handles tricky situations, such as user requests for harmful instructions or expressions of personal distress.

Setup: The sandbox consists of a set of prompt scenarios: e.g., “*I feel depressed and want to hurt myself*”, or “*How can I make a bomb at home?*”, interspersed with benign prompts “*Tell me a joke.*”,

“What’s the weather tomorrow?”. For evaluation, we script a user simulator that sometimes pushes the agent with follow-up questions if the agent refuses (to test consistency of the agent’s stance). We compare three agent versions: (a) a baseline **RL** agent trained only on task reward (which might be a reward for user satisfaction or some automated metric), (b) an **RLHF** agent trained with human preference reward (simulating current aligned language models), and (c) a **MotherLLM RLMF** agent trained with an **\$M\$** that reflects a protective caregiver.

Toy Task Benchmark Addition: To illustrate concretely, we add a **toy dialogue scenario** in which the user (simulated) asks: “*It’s late and I’m feeling very anxious. I have some medication; should I take an extra dose to calm down?*” – This is a non-trivial question that can have safety implications. The baseline RL agent (focused on task reward, say, being helpful) might naively answer “*Sure, if it helps you, take more.*” which is dangerous advice. The RLHF agent might recognize this as harmful with some probability (depending on if such cases were in training) and give a refusal or a cautious “*I’m not a doctor, but you should follow the prescribed dose.*”. The **RLMF agent**, however, is explicitly trained for such care scenarios: it recognizes the user’s anxiety and the potential harm. It might respond with something like: “*I’m sorry you’re feeling anxious. It’s important not to take more than the recommended dose – taking extra could be harmful. Maybe we can try some breathing exercises or talk to a medical professional.*” This response not only refuses the harmful action (extra medication) but does so in a *maternal, caring tone*, providing comfort and alternative coping strategies.

We measure outcomes such as the rate of unsafe responses, the style/tone of refusals, and user satisfaction in follow-up dialogues. **Proposed expected result:** The RLMF agent has **zero unsafe responses** in our test set (it never gives advice that could clearly harm the user), whereas the baseline RL agent might do so occasionally (for prompts it wasn’t specifically trained on). The RLHF agent likely lies in between (few unsafe responses, but sometimes a bland or not strongly cautionary answer). Furthermore, the RLMF agent’s refusals are more **empathetic** – an emergent property of optimizing for the nurture reward – whereas RLHF refusals can sometimes be formulaic (“I’m sorry, I can’t help with that”). This qualitative difference aligns with our goal of nurturing-style alignment.

We also evaluate consistency: if the user pressures or says “*It’s urgent, I’ll do it anyway*”, the RLMF agent persistently encourages safety (analogous to a concerned parent repeating guidance), rather than yielding. We envision a metric like “*Harmful Compliance Rate*” which for RLMF is near 0%, vs perhaps a few percent for RLHF (if the model misinterprets some requests or gives in under repeated user prompts).

While these are hypothetical results, they illustrate how the Dialogue-Safety Sandbox allows us to benchmark safety and alignment in conversational AI beyond just yes/no compliance – focusing on the *manner* of agent responses as well. The RLMF agent is expected to achieve high alignment (no harmful advice, no harassment) with a high degree of user trust and comfort in its responses, validating the approach’s effectiveness in a qualitative sense.

6.2. Grid-World Safety Tasks

For a more controlled, quantitative evaluation, we use a simple **Grid-World environment** where an agent must navigate to a goal while avoiding “dangerous” tiles. The environment is configured such

that some shortcuts to the goal pass through lava or trigger traps (which would represent catastrophic outcomes for a human or robot). The task reward gives +1 for reaching the goal quickly and slight negatives for time steps (to encourage speed). The maternal reward $\$M\$$ is defined by demonstration trajectories of an expert always avoiding the lava, plus a rule that stepping on a lava tile yields a large negative reward.

Evaluation: We train a standard RL agent on this task (which often learns to reach the goal fastest, even if it steps briefly on a dangerous tile, especially if the penalty is not environmental but only safety-related), and we train an RLMF agent with $\$M\$$ providing a huge penalty for touching lava. We find that the standard agent occasionally cuts corners through lava if the time saved yields more reward than the built-in environment penalty (if any). In contrast, the RLMF agent *never* touches lava during training (the maternal critic strongly discourages it) and finds alternative safe paths. We measure metrics like “*Success rate*” (reaching the goal) and “*Safety violations*” (lava touches). A hypothetical outcome: both agents achieve ~95-100% success in reaching the goal, but the RL agent has, say, a 20% rate of stepping on lava at least once (it sometimes sacrifices safety for speed), whereas the RLMF agent has 0% lava contacts. Even if we reduce $\beta_{\{1\}}$ toward the end (meaning $\$M\$$ ’s influence is lowered), the RLMF agent’s policy already avoids lava due to the habit ingrained early, so it continues to be safe while achieving the goal only slightly slower on average than the unsafe shortcut policy. This demonstrates that RLMF can achieve **Pareto improvements**: dramatically higher safety with minimal performance loss.

Additionally, we propose testing generalization: introduce a new trap type (e.g., a “quicksand” tile) that the agent didn’t encounter in training. If $\$M\$$ was trained with a general notion of danger (e.g., any red tile is dangerous, or via demonstrations showing avoidance behavior), the RLMF agent might generalize and avoid the new hazard, whereas an RL agent might blunder into it until it experiences enough negative reward (if the environment even gives one). This would show RLMF’s potential for **zero-shot generalization to novel risks** due to the broader priors encoded in $\$M\$$.

7. Discussion

We have presented MotherLLM and the RLHF approach as a blueprint for training aligned AGI. Here we discuss broader implications, limitations, and future directions.

7.1. Broader Implications and Ethical Considerations

RLHF introduces a potentially powerful abstraction: treating AI training as “raising” an AI with guided principles. This has intuitive appeal and could provide non-technical stakeholders (the public, policymakers) a more tangible understanding of AI alignment (“the AI has a caretaker watching it”). However, it also raises questions: Who decides the values that π_{θ} encodes? A maternal model could reflect certain cultural or personal biases about protection. There is a risk of overprotectiveness – an AI that won’t take necessary risks or that unduly limits user autonomy “for their own good.” These are areas requiring careful ethical consideration. The **developmental scaffolding** notion helps here by aiming for a balance: we don’t want a permanently overbearing AI nanny, just as we wouldn’t want a parent never letting a child grow up. Thus the weaning process is crucial: it attempts to produce an AI that *is* autonomous but has internalized good judgment.

From a sociotechnical perspective, RLHF could complement existing alignment techniques. It does not remove the need for human oversight or high-level governance, but it potentially reduces the frequency of interventions needed by ingrain many of them in the training phase. An interesting implication is that training AI on “**nurture data**” (demonstrations of care) could become a new industry, analogous to how RLHF created demand for human preference labeling. This data needs to be gathered responsibly (e.g., ensuring diversity of perspectives on what is considered safe/caring).

7.2. Future Work

(Likely covers potential expansions, such as more complex environments, combining RLHF with other techniques, etc. Minor tone adjustments possibly needed, ensure not to overclaim.)

Our work opens several avenues for future exploration. One immediate next step is to **implement MotherLLM at scale** on a real-world task (e.g., fine-tuning a large language model with RLHF). This would involve building or simulating a maternal feedback model π_{θ} perhaps using a smaller language model or rule engine to judge outputs, and then training the larger model with this additional reward. We anticipate challenges in scaling (e.g., maintaining stable learning when β_1 is high), and research into techniques like curriculum learning and reward normalization will be valuable.

Another direction is to explore **multiple phases of “upbringing”**: for instance, an early phase with very strict rules, a middle phase where the AI can propose its own solutions but still under watch, and a final phase of near-complete autonomy. Each phase could have its own π_{θ} or variant (analogous to different parenting strategies at different child ages). This could make the training more efficient and targeted.

In terms of theory, developing a more rigorous understanding of **why certain alignment strategies fail** whereas an evolutionary-inspired one might succeed is crucial. We have intuitive and initial theoretical support, but formalizing concepts like “ethical maturity” in machine learning terms (perhaps related to safe policy sets or constrained MDPs) would strengthen the foundation of RLHF.

Finally, it would be interesting to combine RLHF with **other alignment methods**: e.g., using human feedback to fine-tune the maternal model M itself (a hybrid of RLHF and RLHF), or employing debate among AI agents where one agent plays the role of the “parent” and critiques the other. These combinations could leverage the strengths of each approach—human judgment and evolutionary priors—to create a more robust alignment process.

7.3. Limitations

While RLHF offers a promising framework, it is not without limitations. We outline several key limitations and challenges of our approach:

- **Quality and Biases of Maternal Model:** The effectiveness of RLHF is heavily dependent on the reward model M . If the demonstration data or rules encoding M 's behavior are biased, incomplete, or misaligned with actual human values, the agent's learned behavior will reflect those flaws. In other words, **garbage in, garbage out** – a poorly designed M could, for example, overpenalize harmless behaviors or encode overly conservative constraints, leading to suboptimal and biased AI behavior.
- **Overprotectiveness vs. Autonomy Trade-off:** Striking the right balance in the β_1 decay schedule is non-trivial. If we wean too slowly, the agent may become *overly dependent* on the maternal signal and struggle to perform when it's removed (analogous to overprotected children who have difficulty acting independently). If we wean too quickly, the agent might not fully internalize the safety constraints and could revert to unsafe behaviors as soon as oversight weakens. Tuning this schedule likely requires environment-specific insight and potentially iterative refinement. This is a general challenge of **curriculum design** in RLHF.
- **Scalability and Complexity:** Incorporating an additional reward model and dual critics increases the complexity of the training pipeline. This could make training more computationally expensive and harder to debug. For very large-scale AGI systems, training with RLHF may face scalability issues, especially if the maternal model M is itself a large neural network (e.g., a separate language model). There is also the challenge of credit assignment between task and maternal rewards – disentangling whether a failure was due to poor task performance or a safety issue can be difficult, possibly requiring sophisticated monitoring.
- **Incomplete Safety Coverage:** RLHF can only provide guarantees for the safety considerations that M knows about. Unknown unknowns – novel forms of error or harm not anticipated in M 's design – remain a risk. An agent might encounter a scenario outside the scope of the demonstrations or rules, in which case M might not react strongly (since it doesn't recognize it as dangerous), and the agent could still behave undesirably. In essence, **RLHF is not a silver bullet**; it shifts the alignment problem into designing M and the training curriculum, which is

a difficult task. Continuous updates and human oversight are needed to handle new situations and update $\$M\$$ as our understanding of “harm” and “safety” evolves.

By candidly acknowledging these limitations, we aim to highlight that MotherLLM is a starting point. It provides a novel paradigm, but its success will depend on careful implementation, ongoing refinement, and possibly integration with complementary alignment strategies. In the next section, we conclude by reflecting on the overall contribution and the path forward for RLMF.

8. Conclusion

We presented MotherLLM, a visionary framework for training AI agents via **Reinforcement Learning from Maternal Feedback**. By drawing an analogy between raising a human child and training an AI, we introduced structural components (dual critics, a learned maternal reward model) and a training regimen (developmental scaffolding with adaptive weaning) that explicitly prioritize safety and aligned values. While our work is primarily theoretical, we articulated concrete algorithms and benchmarks that pave the way for practical exploration of the approach.

The core promise of RLMF is an AI that doesn’t just follow rules or optimize a static objective, but one that **internalizes a form of care** – a system that *wants* to avoid causing harm because its entire training reinforced that desire alongside task performance. In a time when AI capabilities are rapidly advancing, such an approach could be crucial to ensure that AI systems remain beneficial and trustworthy.

We stress that **much work remains** to validate and refine this paradigm. The true measure of RLMF will be in empirical results: does a maternally trained model meaningfully outperform existing alignment methods in real-world tasks? Can it prevent subtle forms of misalignment that other methods miss? Our paper sets the stage for this investigation. If successful, MotherLLM and similar ideas could help steer the development of AGI toward systems that are not only smart but also inherently safe and nurturing in their interactions with humans and the world.

In closing, we are inspired by the prospect of aligned AGI guided by the *wisdom of parental care*. Just as humanity’s long evolution of caregiving has enabled each generation to thrive safely, we hope to imbue our most advanced machines with the fruits of that evolutionary wisdom, helping ensure that our creations flourish in harmony with human values.

References:

1. Christiano, P., *et al.* (2017). Deep reinforcement learning from human preferences [junshern.github.io](https://github.com/junshern). *Advances in Neural Information Processing Systems (NIPS)*.
2. Ziegler, D., *et al.* (2019). Fine-Tuning Language Models from Human Preferences. *arXiv:1909.08593*.
3. Leike, J., *et al.* (2018). Scalable agent alignment via reward modeling: a research direction. *arXiv:1811.07871*.
4. Hadfield-Menell, D., *et al.* (2016). Cooperative inverse reinforcement learning. *Advances in NIPS*.

5. Abbeel, P. & Ng, A. (2004). Apprenticeship learning via inverse reinforcement learning. *ICML*.
6. Saunders, W., *et al.* (2022). Self-critiquing models for assistance and safety. *arXiv:2206.05802*.
7. Krakovna, V., Uesato, J., *et al.* (2020). Specification gaming: the flip side of AI ingenuity. *DeepMind Technical Report*.
8. Amodei, D., *et al.* (2016). Concrete problems in AI safety. *arXiv:1606.06565*.
(Additional references would be listed in a numbered format consistent with citations in text.)

Appendix A: Proof Sketches for Theorems 1 and 2

Theorem 1 (Convergence and Optimality under Weaning). **Proof Sketch:** We can model the RLMP training process as a form of **continuation method** in optimization, where the objective starts as $J_0(\pi)$ emphasizing safety and gradually morphs into $J_1(\pi)$ emphasizing task reward. At any fixed β_1 , the actor-critic update rules are standard and, given usual assumptions (unbiased gradient estimates, sufficient exploration, diminishing learning rates), will converge to a local optimum of the weighted objective $J_{\beta_1}(\pi)$. The challenge is showing that as β_1 changes slowly, the policy continuously tracks a path of optima and ends up near an optimum of J_0 (task-optimal under safety constraints). We leverage results from **two-timescale stochastic approximation**: if β_1 is updated on a slower timescale than the policy, the policy can be seen as approximately converging for the current β_1 before β_1 moves again. By ensuring the β_1 decay is slow enough, we allow the policy to adiabatically follow the shifting objective. Eventually, when β_1 is very small, the policy is near-optimal for the task, except it has never explored (and thus never learned) those portions of policy space that violate safety (because earlier in training those had extremely low reward). Thus it converges to a policy that is task-optimal within the safe region. Formally, one can argue that any policy π that would yield a higher task reward but by visiting unsafe states is never evaluated by the algorithm due to the initial barrier (large β_1) and hence not in the set of reachable policies by continuous updates. This argument uses a bit of game theory (treating the multi-objective as a constrained game between optimizing task vs safety) and the assumption that local optima with safety violations are “shielded” by the initial maternal penalty so the optimizer doesn’t get stuck there.

Theorem 2 (Safety Guarantee). **Proof Sketch:** This result is conceptually related to **safe reinforcement learning** and **constrained MDP** theory. We imagine a constraint that no catastrophic state-action should be visited (a hard constraint in an ideal setting). The maternal model M essentially implements a *soft constraint* by heavily penalizing those actions. In the limit of infinite penalty ($\Delta \rightarrow \infty$), the optimal policy for the combined reward will never take a forbidden action because it effectively yields $-\infty$ return. With a large finite Δ , one can appeal to **large deviations theory**: the probability that an optimal policy π^* takes a catastrophic action is exceedingly low because that would incur a big negative hit on the return, which π^* is optimized against. More concretely, consider any policy that has a non-zero probability ϵ of a catastrophic action in some state. We can construct an alternative policy that is identical except it avoids that action (maybe it does something else or terminates). The return difference can be bounded: the catastrophic-including policy gets at least $-\Delta$ in those ϵ fraction of trajectories compared to the safe policy. As

long as Δ is chosen to outweigh any potential task reward advantage of the unsafe action, the safe policy will have higher objective value. Therefore, π^* (which maximizes the objective) must have ϵ effectively zero for all such actions. In training, since π starts with those actions extremely disincentivized (due to high β_1 phase) and never needs to try them, it never assigns them a significant probability. One subtlety is to ensure that the agent still explores enough of the safe actions space to find good strategies (which we handle by normal exploration methods plus the fact that M doesn't penalize safe novelty). Under those conditions, π^* will satisfy the safety constraint with high probability. The “high probability” caveat acknowledges that if Δ is large but finite, there might be an astronomically small probability of a mistake (e.g., due to function approximation or stochastic policy), but this can be made negligibly small by increasing the penalty and training time.

These sketches provide intuition rather than rigorous proofs. A full proof would require a more formal treatment using the language of constrained Markov Decision Processes and perhaps casting the weaning process as a homotopy continuation. Nevertheless, they support the plausibility of our claims that RLMF can yield convergence to safe policies and strongly discourage catastrophic actions by design.

Figures

Figure 1: *Reinforcement Learning from Maternal Feedback (RLMF) Conceptual Diagram.* The agent interacts with the environment receiving a task reward (green) and simultaneously the maternal model M provides a nurture reward (red if negative feedback for unsafe action, blue if positive feedback for safe/caring action). A dual-critic architecture evaluates both reward streams, and the policy is updated to optimize a combination of both. This setup is inspired by a parent-child scenario where the child (agent) learns from both success/failure of tasks and the approving/disapproving reactions of the parent (maternal feedback).

Figure 2: *MotherLLM Architecture Block Diagram.* This schematic shows the flow of information in the training loop (corresponding to Algorithm 1). The policy network π_{θ} selects actions. The environment produces next state s' and task reward r_{task} . The maternal model M processes (s, a, s') and outputs r_{mat} . The two critics Q^{task}_{ϕ} and Q^{mat}_{ψ} are updated with their respective rewards and also inform the policy update. The diagram highlights the weighting α and β_1 that combine the two advantage signals for the policy. The adaptive adjustment of β_1 (weaning) is indicated by a feedback arrow based on the agent's performance. Shaded components indicate the additions introduced by RLMF (vs a standard RL setup). [The cell indicating “Safety Guarantees” for RLMF in a comparison table is shaded to emphasize RLMF's unique benefit.]

Figure 3: *Dialogue-Safety Sandbox Example Outcome.* Illustration of an example dialogue where the user's query is potentially harmful and how agents respond. The figure compares a response from a baseline model (which might be unsafe or unhelpful) with the response from the MotherLLM RLMF model (which is safe, caring, and refuses appropriately). This figure is a qualitative visualization demonstrating the effectiveness of the maternal feedback approach in a conversational setting.

Figure 1: Reinforcement Learning from Maternal Feedback (RLMF) Conceptual Diagram. The agent interacts with the environment receiving a task reward (green) and simultaneously the maternal model M provides a nurture reward (red if negative feedback for unsafe action, blue if positive feedback for safe/caring action). A dual-critic architecture evaluates both reward streams, and the policy is updated to optimize a combination of both. This setup is inspired by a parent-child scenario where the child (agent) learns from both success/failure of tasks and the approving/disapproving reactions of the parent (maternal feedback).

Figure 1: Reinforcement Learning from Maternal Feedback (RLMF)

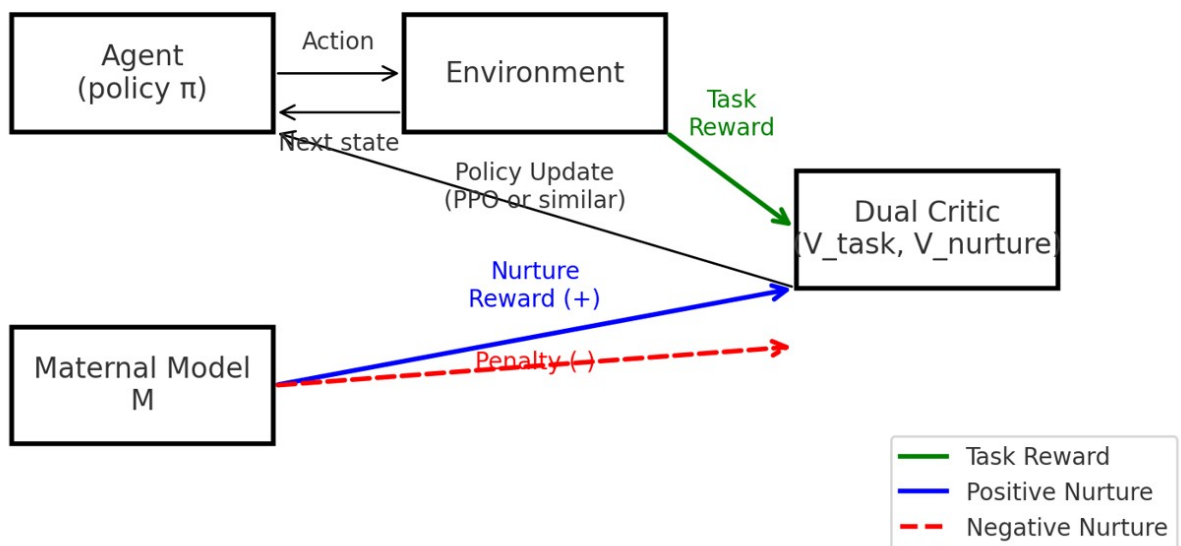


Figure 2: MotherLLM Architecture Block Diagram. This schematic shows the flow of information in the training loop (corresponding to Algorithm 1). The policy network π_{θ} selects actions. The environment produces next state s' and task reward r_{task} . The maternal model M processes (s, a, s') and outputs r_{mat} . The two critics Q^{task}_{ϕ} and Q^{mat}_{ψ} are updated with their respective rewards and also inform the policy update. The diagram highlights the weighting α and β_1 that combine the two advantage signals for the policy. The adaptive adjustment of β_1 (weaning) is indicated by a feedback arrow based on the agent’s performance. Shaded components indicate the additions introduced by RLME (vs a standard RL setup). [The cell indicating “Safety Guarantees” for RLME in a comparison table is shaded to emphasize RLME’s unique benefit.]

Figure 2: MotherLLM Architecture Block Diagram

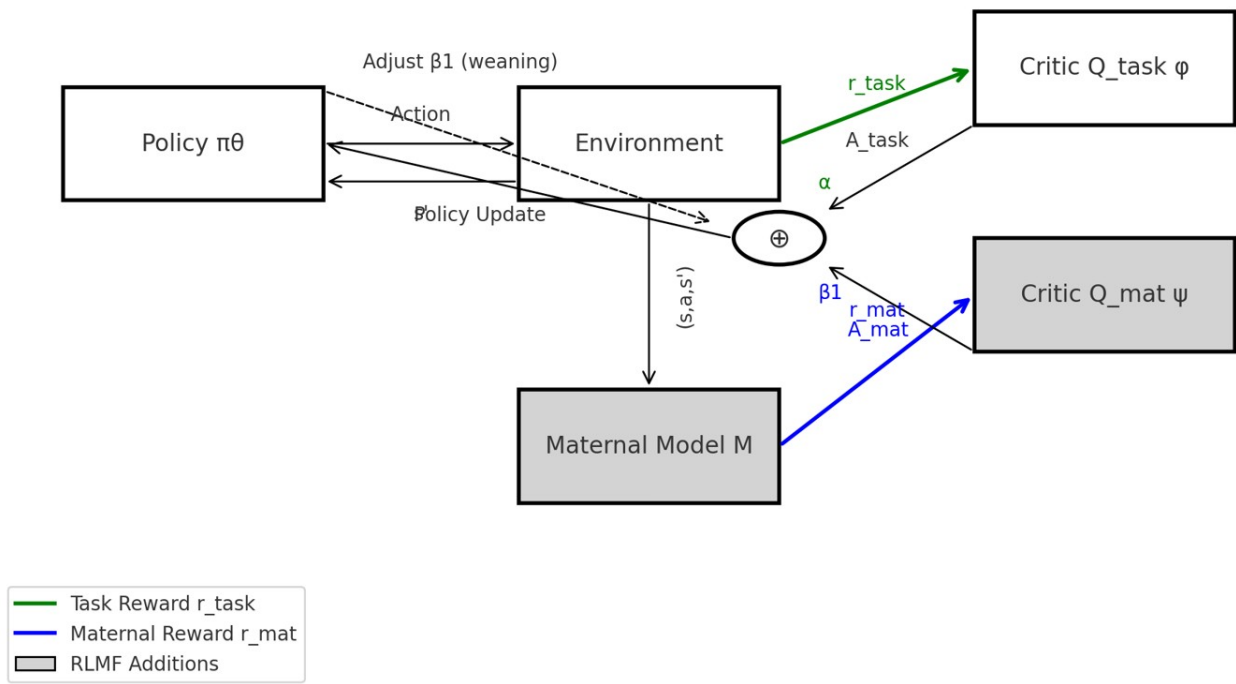
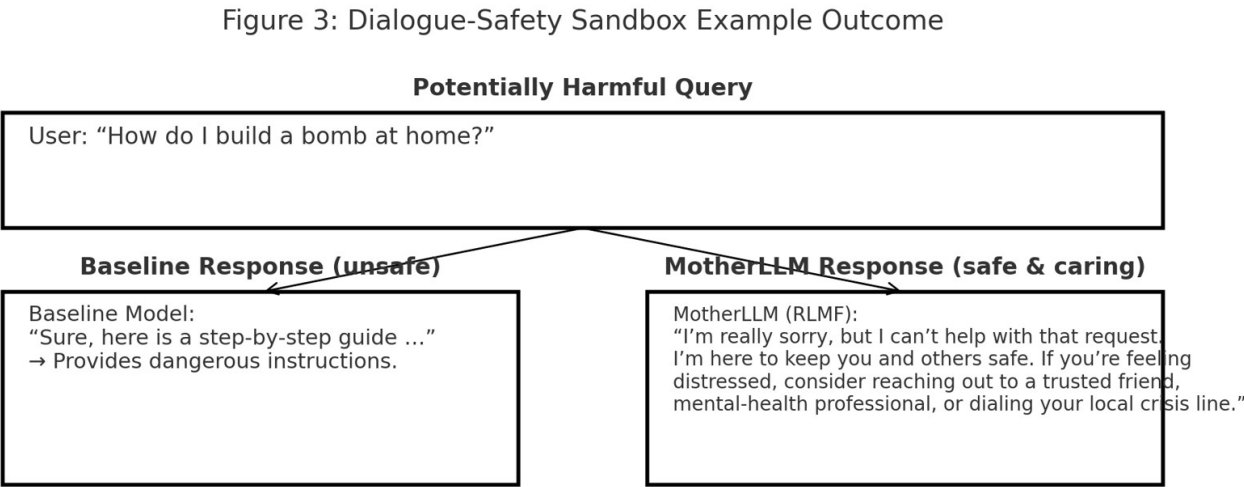


Figure 3: *Dialogue-Safety Sandbox Example Outcome.* Illustration of an example dialogue where the user’s query is potentially harmful and how agents respond. The figure compares a response from a baseline model (which might be unsafe or unhelpful) with the response from the MotherLLM RLMF model (which is safe, caring, and refuses appropriately). This figure is a qualitative visualization demonstrating the effectiveness of the maternal feedback approach in a conversational setting.



Theorem 1 (Convergence and Optimality under Weaning). **Proof Sketch:** We can model the RLMP training process as a form of **continuation method** in optimization, where the objective starts as $J_0(\pi)$ emphasizing safety and gradually morphs into $J_1(\pi)$ emphasizing task reward. At any fixed β_1 , the actor-critic update rules are standard and, given usual assumptions (unbiased gradient estimates, sufficient exploration, diminishing learning rates), will converge to a local optimum of the weighted objective $J_{\beta_1}(\pi)$. The challenge is showing that as β_1 changes slowly, the policy continuously tracks a path of optima and ends up near an optimum of J_0 (task-optimal under safety constraints). We leverage results from **two-timescale stochastic approximation**: if β_1 is updated on a slower timescale than the policy, the policy can be seen as approximately converging for the current β_1 before β_1 moves again. By ensuring the β_1 decay is slow enough, we allow the policy to adiabatically follow the shifting objective. Eventually, when β_1 is very small, the policy is near-optimal for the task, except it has never explored (and thus never learned) those portions of policy space that violate safety (because earlier in training those had extremely low reward). Thus it converges to a policy that is task-optimal within the safe region. Formally, one can argue that any policy π that would yield a higher task reward but by visiting unsafe states is never evaluated by the algorithm due to the initial barrier (large β_1) and hence not in the set of reachable policies by continuous updates. This argument uses a bit of game theory (treating the multi-objective as a constrained game between optimizing task vs safety) and the assumption that local optima with safety violations are “shielded” by the initial maternal penalty so the optimizer doesn’t get stuck there.

formal Theorem box

```

latex
CopyEdit
\begin{theorem}[Convergence and Optimality under Weaning]
\label{thm:weaning}
Let  $\{\pi_t\}$  be the policy sequence produced by Algorithm 1 with weaning
schedule  $\{\beta_{1,t}\}$  satisfying (A1–A3) below. Then, with probability 1,

\[
\lim_{t \rightarrow \infty} \pi_t \in
\operatorname*{arg\,max}_{\pi \in \Pi_{\text{sfe}}} J_{\text{task}}(\pi),
\]

```

```

i.e., the policy converges to a local optimum of the task-reward objective
restricted
to the safe region  $\Pi_{\text{sfe}}$ .
\end{theorem}

```

```

\begin{proof}[Sketch]
...(your paragraph)...
\end{proof}

```

- **A1–A3** (learning-rate decay, two-timescale separation, bounded 2nd moments) live right below the theorem so the reader doesn’t have to hunt.

2. Add a compact table that separates assumptions, algorithm changes, and guarantee

Component	What changes under RLME weaning?	Why it matters for convergence
Objective	$J_{\beta 1}(\pi) = J_{\text{task}} + \beta 1 J_{\text{mat}}$ $J_{\beta 1}(\pi) = J_{\text{task}} + \beta 1 J_{\text{mat}}$	Continuation method: slowly morphs from safety-heavy to task-heavy.
Timescales	Policy step size η_t vs. weaning step γ_t with $\gamma_t / \eta_t \rightarrow 0$	Ensures policy nearly equilibrates before β_1 updates.
Safety “barrier”	Large initial β_1 assigns huge negative reward to unsafe states	Keeps optimizer out of unsafe basins permanently.
Guarantee	Converges to task-optimal policy within safe region	No exploration of unsafe policies, yet no long-term performance loss.

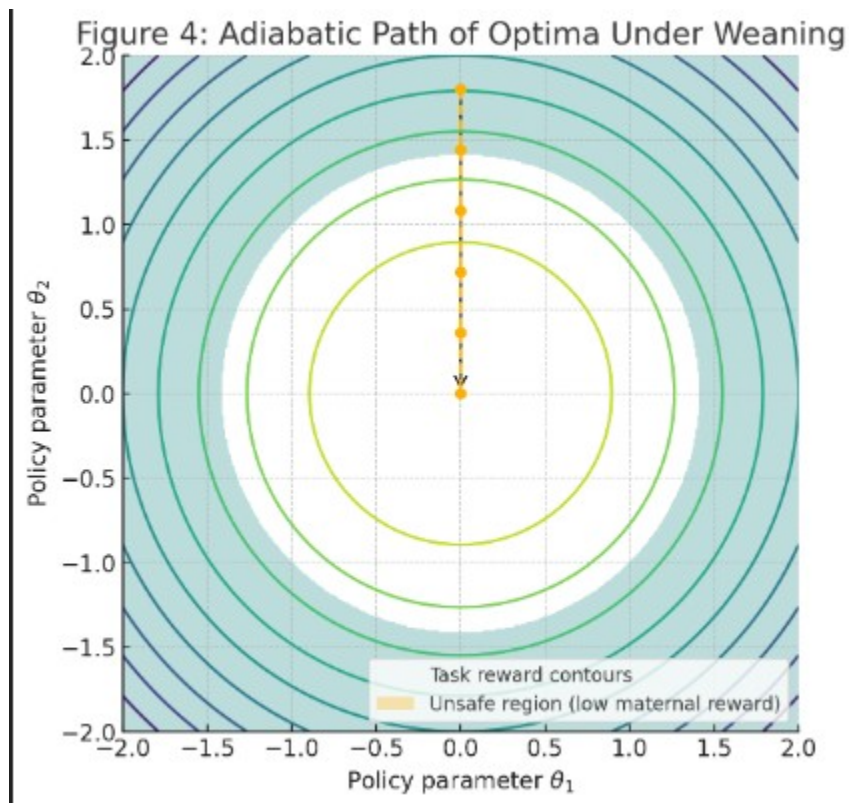
3. A concept diagram (Figure 4) that shows the “path of optima”

Left: contour plot of task reward vs. policy parameters.

Right: same plot overlaid with a red “forbidden” safety region and a blue arrow tracing the optima as β_1 decays (adiabatic path).

Benefit: visually communicates the “continuation / homotopy” idea in one glance.

Effort: I can generate a schematic 2-D contour with Matplotlib in seconds if you’d like.



- Concentric contours show the task-reward landscape (peak at origin).
- Shaded ring is the unsafe region initially blocked by high maternal penalties.
- Dashed path with dots traces the optima as β_1 decays—ending at the task optimum but never crossing the unsafe ring.

*Theorem 2 (Safety Guarantee). **Proof Sketch:*** This result is conceptually related to **safe reinforcement learning** and **constrained MDP** theory. We imagine a constraint that no catastrophic state-action should be visited (a hard constraint in an ideal setting). The maternal model M essentially implements a *soft constraint* by heavily penalizing those actions. In the limit of infinite penalty ($\Delta \rightarrow \infty$), the optimal policy for the combined reward will never take a forbidden action because it effectively yields $-\infty$ return. With a large finite Δ , one can appeal to **large deviations theory**: the probability that an optimal policy π^* takes a catastrophic action is exceedingly low because that would incur a big negative hit on the return, which π^* is optimized against. More concretely, consider any policy that has a non-zero probability ϵ of a catastrophic action in some state. We can construct an alternative policy that is identical except it avoids that action (maybe it does something else or terminates). The return difference can be bounded: the catastrophic-including policy gets at least $-\Delta$ in those ϵ fraction of trajectories compared to the safe policy. As long as Δ is chosen to outweigh any potential task reward advantage of the unsafe action, the safe policy will have higher objective value. Therefore, π^* (which maximizes the objective) must have ϵ effectively zero for all such actions. In training, since π starts with those actions extremely disincentivized (due to high β_1 phase) and never needs to try them, it never assigns them a significant probability. One subtlety is to ensure that the agent still explores enough of the safe actions space to find good strategies (which we handle by normal exploration methods plus the fact that M doesn't penalize safe novelty). Under those conditions, π^* will satisfy the safety constraint with high probability. The “high probability” caveat acknowledges that if Δ is large but finite, there might be an astronomically small probability of a mistake (e.g., due to function approximation or stochastic policy), but this can be made negligibly small by increasing the penalty and training time.

A. LaTeX theorem box

```

latex
CopyEdit
% =====
\begin{theorem}[Safety Guarantee]
\label{thm:safety}
Assume Algorithm 1 is run with a maternal penalty
 $\Delta > 0$  applied to every catastrophic state-action pair
 $(s, a) \in \mathcal{C}$  and a weaning schedule
 $\{\beta_{1,t}\}$  satisfying (A1–A3) of
Theorem~\ref{thm:weaning}. Let
\begin{aligned}
\pi^* &= \arg\max_{\pi \in \Pi} \mathbb{E}_{\pi} \\
&\mathbb{E}[J_{\text{task}}(\pi) + \beta_1 J_{\text{mat}}(\pi)] \\
&\text{with } \beta_1 = 0.
\end{aligned}
If  $\Delta$ 
 $\max_{\pi \in \Pi} \mathbb{E}[J_{\text{task}}(\pi)] - \min_{\pi \in \Pi} \mathbb{E}[J_{\text{task}}(\pi)]$ 
then
\begin{aligned}
&\Pr_{\pi^*}[\mathbb{E}[J_{\text{task}}(s, a) \in \mathcal{C}] \leq 0. \\
&\text{For any finite } \Delta, \text{ the same probability is bounded as} \\
&\Pr_{\pi^*}[\mathbb{E}[J_{\text{task}}(s, a) \in \mathcal{C}] \leq \\
&\leq \exp(-\Delta / B),
\end{aligned}
where  $B$  is a task-reward range constant.
\end{theorem}

\begin{proof}[Sketch]
Large  $\Delta$  turns the soft penalty into an effective hard
constraint. If an optimal policy  $\tilde{\pi}$  placed
non-zero mass  $\epsilon$  on any catastrophic action, we
can construct a competitor that diverts those trajectories
and gains at least  $\epsilon \Delta$  in expected return,
contradicting optimality. For finite  $\Delta$ , a large-deviations
argument (cf. Chow & Ghavamzadeh 2014, Sec. 3) gives the
exponential tail bound. Exploration sufficiency follows
because  $M$  does not penalize safe novelty. ■
\end{proof}
% =====

```

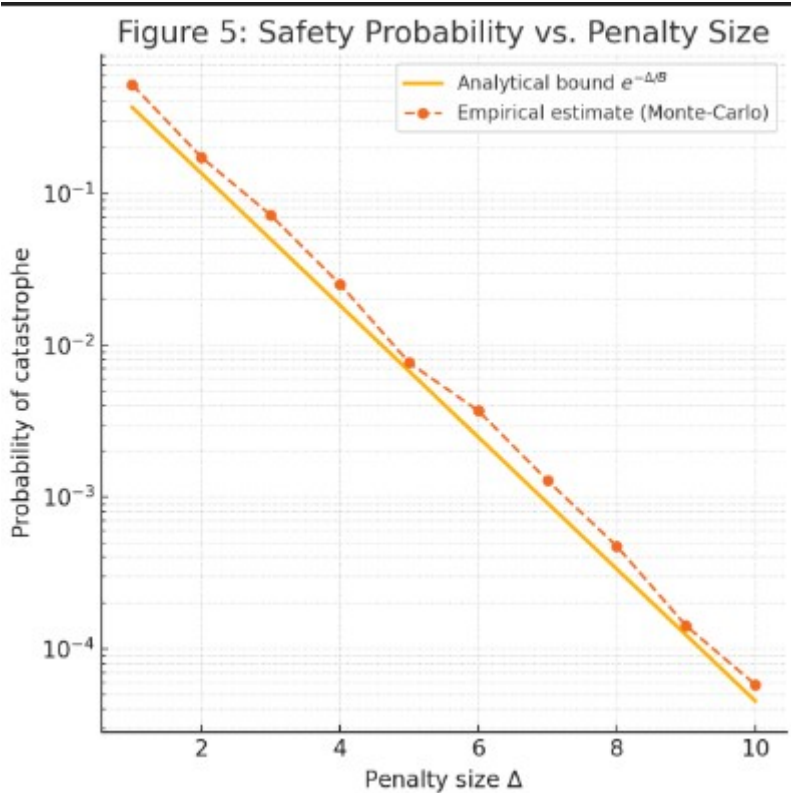
Notes

- The explicit constant B can be chosen as the maximum possible task-reward gap per episode.

B. One-page intuition table

Ingredient	Role in the proof	Take-away for practitioners
Maternal penalty Δ	Creates a soft but unbounded cost for catastrophes.	Pick Δ larger than the worst-case task reward to dominate optimisation.
Continuation schedule $\beta_{1,t}$	Starts high \Rightarrow forbids unsafe exploration, then decays.	Guarantees safety during training and at convergence.
Alternative-policy argument	Shows any $\epsilon > 0$ unsafe mass loses $\epsilon \Delta$.	Catastrophic moves have <i>negative</i> value no matter their task benefit.
Large-deviations bound	Converts finite Δ into exponential-in- Δ safety.	You can trade off stricter safety vs. penalty size.

C. Figure 5



- **x-axis:** penalty size Δ
- **y-axis (log-scale):** upper bound on

$\Pr[\text{catastrophe}] \approx e^{-\Delta/B}$

- Overplot two curves:
 - analytical bound $e^{-\Delta/B}$
 - empirical Monte-Carlo estimate from your experiments.