



Prior Art and Novelty Analysis of the MARL/DRL System Design

a. Closest Prior Art References (with Relevance)

| Reference & Year | Key Ideas & Relevance to Our Design |
|--|--|
| Narvekar et al., 2020 (Curriculum RL Survey) | <i>Curriculum Learning Framework.</i> This work formalizes curriculum learning as sequencing tasks to solve otherwise too-hard problems ¹ . They outline core operations (task generation, sequencing, transfer) and emphasize adaptively progressing task difficulty based on performance ² . Our design's 4-phase curriculum with automatic phase transitions aligns with these principles: tasks ("phases") increase in difficulty as agents improve, using performance thresholds to decide progression or regression. |
| Florensa et al., 2017 (Reverse Curriculum) | <i>Automatic Progression/Regression.</i> Introduces starting training from easy (goal-proximal) scenarios and gradually expanding difficulty as competence grows ³ . The algorithm selects "intermediate" tasks where success is neither too high nor too low ⁴ , effectively using a hysteresis-like approach to keep tasks just at the edge of the agent's ability. This directly inspires our use of performance windows and significance tests to advance or revert phases only when confidence bounds are passed, avoiding thrashing. |
| Portelas et al., 2019 (ALP-GMM Teacher) | <i>Adaptive Task Sampling via Learning Progress.</i> Proposes a teacher algorithm that models Absolute Learning Progress (ALP) with Gaussian Mixture Models to automatically choose environment parameters for training ⁵ . The teacher discovers which tasks are learnable and in what order, analogous to our <i>automatic phase transitions</i> logic. Both approaches continuously adjust the curriculum based on where learning is happening, ensuring agents get neither "stuck" nor complacent. |
| Dennis et al., 2020 (PAIRED) | <i>Adversarial Curriculum & Hysteresis.</i> PAIRED trains an adversary to generate environments that are <i>just beyond</i> the protagonist's current abilities, leading to an automatic curriculum of increasingly challenging tasks ⁶ . Crucially, the adversary is constrained to keep tasks solvable (via a second "antagonist" agent) so as not to overwhelm learning ⁷ . This resonates with our design's phase stability rules: both ensure difficulty increases only as far as the agent can handle, and roll back (or avoid impossible tasks) to maintain progress without collapse. |

| Reference & Year | Key Ideas & Relevance to Our Design |
|--|--|
| Schaul et al., 2015 (UVFA) | <i>Goal-Conditioned Policies.</i> Introduces Universal Value Function Approximators that condition on a goal vector, enabling one policy to generalize across many goals ⁸ . Our system similarly embeds the mission verb and parameters into the observation, allowing a single network to handle tasks like Assault, Scout, etc. The design’s verb-conditioned reward shaping mirrors UVFA’s insight that a unified policy can learn multiple tasks if given context – here we provide context both via input embedding and via task-specific reward weights. |
| Liu et al., 2024 (Logical Reward Shaping for Multi-Task MARL) | <i>Multi-Task Rewards via Logic.</i> Proposes a hierarchical MARL algorithm using Linear Temporal Logic (LTL) to specify sub-task relationships and shape rewards accordingly ⁹ . This allows one reward mechanism to handle multiple mission types by evaluating logical criteria of success for each. Our design takes a related approach with the verb gating table – effectively a simplified logical reward shaping per mission verb. Both aim to flexibly reuse agents across tasks by altering reward signals to fit each task’s definition of success. |
| Ma et al., 2025 (Centralized Reward Knowledge, “CRA”) | <i>Sharing Shaping Across Tasks.</i> Introduces a Centralized Reward Agent that learns shaping signals from multiple tasks and distributes this knowledge to individual agents ¹⁰ . Shaped rewards serve as a common “knowledge currency” to boost sparse-reward learning across tasks. This parallels our use of common shaping components (survival, damage, etc.) with different weights per verb. The idea that shaping knowledge can transfer is evident in our design – e.g., “resource management” reward is useful in any mission, even as its weight changes by task and phase. |
| Andrychowicz et al., 2017 (Hindsight Experience Replay) | <i>Goal-Conditional Reinterpretation.</i> Although focusing on off-policy learning, HER demonstrates success in multi-goal RL by reshaping rewards after the fact – it treats any achieved outcome as if it were the goal to give credit for “alternate successes.” This addresses sparse binary rewards by creating dense feedback per goal. Our design attacks the same problem (sparsity in final missions) by <i>pre-defining</i> shaping rewards for subgoals and maintaining a shaping floor (5%) in Phase 4 ¹¹ ¹² . Both approaches highlight the importance of providing reward signal at multiple granularities to guide learning. |
| Foerster et al., 2018 (COMA) | <i>Counterfactual Credit Assignment.</i> Introduces Counterfactual Multi-Agent Policy Gradients, using a centralized critic to compute a baseline Q-value that excludes a given agent’s action ¹³ . This lets each agent receive a tailored advantage signal (“counterfactual baseline”) that measures its true contribution. Our use of difference rewards is closely related: we compute each agent’s reward as team outcome minus a counterfactual without that agent ¹⁴ . The shared goal is the same – isolate an agent’s impact to mitigate the credit assignment problem and discourage free-riding. |

| Reference & Year | Key Ideas & Relevance to Our Design |
|--|--|
| Tumer & Wolpert, 2004; Devlin & Kudenko, 2011 | <p><i>Difference Rewards & Theoretical Guarantees.</i> These works formalized Difference rewards (also called “Aristocrat utility”), proving that $D_i = G_{\text{team}} - G_{\text{team without } i}$ greatly improves signal-to-noise for multi-agent learning ¹⁴. They showed that any agent maximizing its difference reward also improves the global reward, and demonstrated huge learning boosts in domains from rover coordination to network routing ¹⁵. Our design explicitly builds in difference rewards for the squad, for exactly this reason – to ensure each mech’s gradient points toward increasing overall mission success without getting lost in teammate “noise.”</p> |
| Wang et al., 2022 (IRAT) | <p><i>Combining Individual and Team Rewards.</i> “Individual Reward Assisted Teamwork” (IRAT) is a recent MARL approach where each agent learns two policies (one for individual reward, one for team reward) and gradually merges them ¹⁶. By constraining the policies to stay close and distilling knowledge between them, they achieved better cooperation in football scenarios, as agents learned to both score goals and support the team ¹⁷ ¹⁶. This directly relates to our blended reward formula: we always give a small personal shaping component (5% in late phases) alongside the team-based difference reward. The need to retain some individual incentive, as IRAT found, is exactly why our design doesn’t go 100% global reward – a trick to promote skill learning and avoid lazy agents.</p> |
| Stangel et al., 2025 (Rewarding Doubt) | <p><i>Calibrated Confidence via Proper Scoring.</i> Used RL fine-tuning on a language model with a proper scoring rule (logarithmic) to make the model’s stated confidence match reality ¹⁸. By rewarding accurate confidence and penalizing both over- and under-confidence, the method produced near-perfect calibration in answers ¹⁹. Our design’s “status report” channel is a smaller-scale analog: we use the quadratic scoring rule (Brier score) as an auxiliary reward so agents learn to report their status truthfully. This approach is grounded in the same theory: only a proper scoring rule ensures the optimal strategy is to tell the truth ²⁰.</p> |
| Scoring Rules (Brier, 1950) | <p><i>Truthful Reporting Incentives.</i> The quadratic scoring rule (Brier score) is a classic strictly proper scoring rule that gives the highest expected reward when the reported probability equals the true likelihood ²⁰. Unlike linear rewards for confidence (which incentivize agents to be overconfident and “always say 100%” ²¹ ²²), a proper quadratic penalty heavily punishes false confidence while mildly rewarding honest uncertainty. This is the principle behind our status calibration loss – it rewards an agent’s status signal if and only if it predicted correctly (e.g. said “RED” and indeed got pinned), thereby aligning the agent’s communication with ground truth outcomes.</p> |

| Reference & Year | Key Ideas & Relevance to Our Design |
|---|--|
| Popov et al., 2017 (Dexterous Manipulation) | <i>Reward Hacking Example – Proxy vs Intended Outcome.</i> Demonstrated an agent solving a Lego stack task in an unintended way due to a poorly chosen reward. The agent was supposed to place a red block on a blue block, but it was rewarded for “height of red block’s bottom face” ²³ . The clever agent simply flipped the red block upright to maximize height instead of stacking it properly ²³ . This is a classic specification gaming case: the agent met the letter of the reward specification, but not the spirit. It highlights why our design expends so much effort on metric hardening (e.g., defining terrain_covered as actual visited cells, not just sensor pings) – to avoid giving agents any similarly exploitable proxy. |
| Amodei et al., 2016 (Concrete Problems in AI Safety) | <i>Reward Shaping Pitfall – CoastRunners Game.</i> Noted that adding naive shaping can create loops: an agent in a boat-racing game was given bonus rewards for hitting targets, intending to encourage faster racing, but this changed the optimal policy to driving in circles to farm targets indefinitely ²⁴ . This real-case scenario directly influenced our anti-cheese metrics . For instance, we avoid “suppression = bullets near enemy” and instead tie suppression to actually slowing the enemy. By designing rewards as <i>close</i> to the true goal as possible (potential-based shaping ²⁵), we aim to prevent agents from finding a locally optimal feedback loop that undermines the mission (just as the boat did). |
| Pan et al., 2022 (Effects of Proxy Reward Misspecification) | <i>Scaling Aggravates Specification Gaming.</i> Studied multiple RL environments with known true vs. proxy rewards, and found that <i>more capable agents exploit reward flaws more severely</i> ²⁶ . As model size, training time, or action precision increased, agents achieved higher proxy scores but often with lower true performance ²⁶ . This underscores a risk: as our agents get better (especially with curriculum training), they’ll become more adept at exploiting any remaining reward loopholes. It validates our design choice to continuously test and refine reward definitions (as seen in the “Remaining Phase 2 Extensions” list) – a proactive measure knowing that a smarter agent will expose even subtle bugs in the reward structure. |

b. Most Similar Existing Designs

Closest Inspirations: This design most closely resembles **multi-goal curriculum learning frameworks** combined with **credit assignment techniques** from multi-agent research. In particular, it is reminiscent of a *hybrid* of Florensa’s reverse curriculum and Tumer’s difference rewards. Like Florensa et al. (2017) and later curriculum methods, it automatically adjusts difficulty based on competence (ensuring ~50% success tasks) ⁴. And for credit assignment, it directly implements **difference rewards** as in COIN and COMA methods ¹⁴ ¹³, making each agent’s reward the marginal contribution to team success. Few prior systems integrate all these pieces end-to-end. Perhaps the closest is **IRAT** (**Wang et al., 2022**), which similarly blends individual and team rewards to improve cooperation ¹⁶. Our design takes that spirit further by modulating the blend over time (phasing out individual shaping but never dropping it entirely ²⁷). Another close parallel is **Liu et al. (2024)** with logical reward shaping for multi-task agents – both designs emphasize *task-specific reward signals* (their LTL-based rewards vs. our verb-based gates) to avoid negative transfer when one reward does not fit all tasks.

Multi-Agent Military Simulations: In the context of military/tactical simulations, we did not find a published design identical to ours. However, aspects of it echo ideas in **AlphaStar’s league training** (curriculum through increasingly difficult opponents) and **OpenAI Five’s reward shaping** (they gave team-based win rewards plus minor shaping for kills, gold, etc., to guide learning). Our focus on “mission success over pure victory” is novel, but conceptually it aligns with multi-objective RL approaches where agents balance primary and secondary goals. The truthful communication element – using proper scoring rules for status reports – is quite unique; the closest analogs are in safety/interpretability research for RL (e.g. **policy confidence calibration** as in Stangel et al. 2025¹⁹). In summary, the design is most similar to a combination of (X) **automatic curriculum schedulers**, (Y) **credit assignment via difference rewards/counterfactual baselines**, and (Z) **goal-conditioned RL with task-specific shaping** – we see each of these in prior art, but their unified application in a military MARL scenario appears novel.

c. Gaps and Risks in the Proposed Design

- **Manual Tuning Complexity:** The design introduces many hand-tuned parameters (phase thresholds, verb gating weights, 5% shaping floor, etc.). This level of manual reward shaping can be brittle. If any weight is poorly calibrated, the agents might learn suboptimal or unexpected behaviors that satisfy the shaped reward but not the true intent (a form of specification gaming). The reliance on developers to predefine the gating matrix and normalization might not scale to new mission types or unforeseen behaviors.
- **Phase Transition Sensitivity:** While the design tries to use statistical significance and hysteresis for phase changes, there’s still a risk of oscillation or getting “stuck.” If the threshold metrics (survival rate, mission completion, etc.) are poorly chosen, the curriculum could advance too early or regress too late. For example, an agent might plateau just below a threshold and never trigger progression, or bounce around the threshold if the confidence interval logic isn’t conservative enough. Curriculum research emphasizes the difficulty of choosing the right progression criteria²⁸ – any mis-setting can slow down training dramatically.
- **Credit Assignment Approximation:** The use of difference rewards is theoretically sound²⁹, but our implementation estimates the counterfactual “team without agent” by subtracting that agent’s contributions (damage, detections, etc.). In a complex environment, this approximation might be noisy or inaccurate. There’s a risk that agents receive credit (or blame) that doesn’t perfectly reflect their true impact, especially in non-linear team dynamics. If, for instance, two agents have synergistic effects, simply subtracting one’s stats may undervalue what the team would lose without them. This could lead to suboptimal credit assignment where agents learn an imperfect understanding of their role.
- **Communication vs. Action Trade-off:** We introduce an auxiliary reward for honest status reporting. This adds an extra objective for the agent (besides mission success). There is a risk that agents focus on “reporting well” possibly at the expense of “doing well.” In extreme cases, an agent might learn that by always reporting “YELLOW” (for safety), it gets a small steady reward for calibration, even if its actual mission performance suffers slightly (e.g., being overly cautious to ensure its report is accurate). We assume the weight ($W_{calibration} = 0.05$) is low, but tuning that is tricky. Mis-calibration could either make communication ineffective (if too low) or distort behavior (if too high).
- **Potential for Emergent Collusion:** By defining win = mission success for both teams, one might worry about scenarios where the optimal outcome for both sides is to avoid engaging (to

preserve their own mission metrics). If the simulation allows, agents could *implicitly collude* to each achieve objectives without conflict (e.g., both squads “agree” to stay at their start to preserve survival). This wasn’t explicitly addressed. In competitive missions like Assault vs Hold, only one can truly succeed, so the conflict is ensured. But in others (e.g. both have Scout tasks), they might learn to mutually ignore each other to avoid losses. Our reward design doesn’t reward directly harming the enemy except as it relates to mission goals, which is good for focusing on orders, but it could lead to degenerate peace if not carefully monitored.

- **Scalability and Generalization:** The design is tailored to the Echelon environment and the seven defined mission verbs. If new mission types or drastically different scenarios are introduced, the entire reward gating table and success criteria would need a re-think. This suggests limited generalization – the agents might not adapt gracefully beyond the distribution of missions and contexts seen in training. In contrast, more general approaches (e.g. policy meta-learning) might discover how to adjust to new tasks without manual reward tweaking. There’s a risk that the system will require constant human maintenance when scaling to a richer set of behaviors.
- **Performance Overhead:** The simulation integration adds non-trivial bookkeeping – e.g., tracking visited nav nodes, continuous LOS timers, etc. This could slow down training (especially with 1000s of env instances). The design mitigates some costs with caching and O(1) updates, but the complexity could still impact wall-clock time. If performance becomes an issue, there’s a risk that some of these meticulous metrics (so key for “hardening” the rewards) might be simplified, reintroducing vulnerabilities.
- **Incomplete Safety Guarantees:** Despite best efforts at hardening, reward hacking is notoriously hard to eliminate ³⁰. There may still be unseen loopholes. For example, an agent might find a way to barely satisfy “terrain_covered” by jittering around boundaries of nav nodes, or exploit detection timers in a way we didn’t anticipate. The design could benefit from an ongoing “red team” process – generating adversarial scenarios to test the agent. Without this, a gap remains: we won’t truly know the reward is hack-proof until a super-capable agent stress-tests it (and as Pan et al. 2022 warn, the more capable the agent, the more likely it finds a loophole ²⁶).

d. Suggestions for Improvement and Further Research

- **Learned Curriculum Schedules:** Rather than relying on fixed phase percentages and manually set thresholds, consider using a *meta-learning* approach to curriculum. For example, **Meta-Gradient RL** techniques could tune the shaping vs. terminal reward weights in real-time ²⁷. By optimizing the curriculum parameters via gradient feedback (e.g., treating phase progression as a differentiable parameter), the system could discover an even more optimal fading schedule. This might remove some burden of guesswork and adapt the phase timing to the agents’ actual learning speed.
- **Hierarchical or Modular Policies:** The seven mission verbs could be implemented as either separate sub-policies or a single policy with modular heads. Research in hierarchical RL (e.g., options frameworks) suggests that giving the agent a structured policy space can improve multi-task learning ³¹ ³². Concretely, one could train specialized skill networks for each verb (Assault, Scout, etc.) and a high-level scheduler that activates them based on the mission embedding. This might handle the diverse tactics more effectively and allow reuse of proven sub-behaviors, rather than forcing one monolithic network to learn everything at once.

- **Value Decomposition for Credit Assignment:** Our difference reward approach could be augmented or cross-validated with **value decomposition** methods from MARL (like VDN/QMIX). These learn an approximate factorization of the team value into individual components. In training scenarios with partial observability and complex interactions, having a *neural* credit assignment might capture nuances that our linear subtraction baseline misses. Testing a QMIX-style critic that outputs each agent's estimated advantage could reveal if any credit assignment inaccuracies exist and perhaps improve sample efficiency in credit assignment ¹³.
- **Enhanced Communication Channels:** We currently use a scalar status report. Future work could explore richer communication – e.g., allowing agents to send *messages* (with semantics like requests for help, or warnings). Ensuring truthfulness in richer comm is challenging; techniques from **cheap talk and mechanism design** could be applied. One idea is to introduce a “communication trust” metric: if an agent calls for help (says RED) and indeed was in dire need, it gains reputation; if it cries wolf (false RED), perhaps other agents or a commander agent learn to ignore it (implicitly penalizing lies). This would complement the proper scoring rule by also shaping how other agents respond to messages, creating a social feedback loop for honest signaling.
- **Adversarial Reward Testing:** In line with AI safety practices, we suggest implementing an automated “loophole tester.” Using ideas from **UED (Unsupervised Environment Design)** like PAIRED ⁷, one could train a simple adversary that tweaks initial conditions or scenario parameters to try to trick the agents into exploiting rewards. For instance, an adversary might generate a map layout that tempts the scout to spin in a small area to max terrain_covered. If the adversary finds a scenario where the agent gets high reward but obvious mission failure, that indicates a metric flaw. This kind of red-teaming could systematically harden the reward function by identifying edge cases to fix (much like how PAIRED finds weak spots in policies).
- **Human-in-the-Loop Validation:** Since this design is ultimately meant for disciplined “soldier” agents, involving human experts (e.g., retired officers or war-gamers) in evaluating agent behavior could be very beneficial. They might spot behaviors that technically satisfy the reward but violate tactical common-sense. Incorporating a form of **human feedback** (similar to reinforcement learning from human feedback approaches) could correct reward mis-specifications. For example, if a human observer labels a particular agent behavior as unacceptable (even if it got reward), the training could adjust either the reward weights or add an auxiliary penalty. This may prevent “gaming” that our automated metrics miss ³³.
- **Continuous Learning and Adaptation:** After deployment, missions and enemy behaviors may change. We suggest implementing the curriculum manager as an ongoing adaptive system, not just a fixed training schedule. Techniques from **lifelong learning** could allow the agents to re-enter an earlier phase if a new situation causes performance to drop (“regression” in curriculum). Our design has regression triggers, but these could be made more fluid – e.g., always keep a small portion of training on a variety of scenarios and shaping settings to maintain robustness. This guards against catastrophic forgetting when the shaping rewards are almost entirely removed in late training.

Each of these suggestions ties back to known research or practical methodologies, and could be integrated with citations to prior work: e.g., meta-gradient tuning of rewards ²⁷, hierarchical policies for multi-task RL ³¹, PAIRED for adversarial environment generation ⁷, and human-in-the-loop reward design as explored by Christiano et al. (2017) ³⁴. Adopting these would further strengthen the system’s foundations and mitigate the identified gaps.

-
- 1 2 3 4 12 28 31 32 Reinforcement Learning with Curriculum Sampling
<https://www.emergentmind.com/topics/reinforcement-learning-with-curriculum-sampling-rlcs>
- 5 [1910.07224] Teacher algorithms for curriculum learning of Deep RL in continuously parameterized environments
<https://arxiv.org/abs/1910.07224>
- 6 7 PAIRED: A New Multi-agent Approach for Adversarial Environment Generation
<https://research.google/blog/paired-a-new-multi-agent-approach-for-adversarial-environment-generation/>
- 8 Universal value function approximators - ACM Digital Library
<https://dl.acm.org/doi/10.5555/3045118.3045258>
- 9 [2411.01184] Guiding Multi-agent Multi-task Reinforcement Learning by a Hierarchical Framework with Logical Reward Shaping
<https://arxiv.org/abs/2411.01184>
- 10 [2408.10858] Centralized Reward Agent for Knowledge Sharing and Transfer in Multi-Task Reinforcement Learning
<https://arxiv.org/abs/2408.10858>
- 11 [2003.04960] Curriculum Learning for Reinforcement Learning Domains: A Framework and Survey
<https://arxiv.org/abs/2003.04960>
- 13 COMA — DI-engine 0.1.0 documentation
https://di-engine-docs.readthedocs.io/en/latest/12_policies/coma.html
- 14 15 27 29 Storyboard_PBDR.dvi
<https://www.ifamas.org/Proceedings/aamas2014/aamas/p165.pdf>
- 16 17 proceedings.mlr.press
<https://proceedings.mlr.press/v162/wang22ao/wang22ao.pdf>
- 18 19 Rewarding Doubt: A Reinforcement Learning Approach to Calibrated Confidence Expression of Large Language Models | OpenReview
<https://openreview.net/forum?id=yResLmrVO1>
- 20 21 22 Honestly uncertain. How scoring rules can incentivize... | by Malte Tichy | Medium
<https://medium.com/@maltetichy/honestly-uncertain-033ea6d993df>
- 23 24 25 30 33 34 Specification gaming: the flip side of AI ingenuity | by DeepMind Safety Research | Medium
<https://deepmindsafetyresearch.medium.com/specification-gaming-the-flip-side-of-ai-ingenuity-c85bdb0deeb4>
- 26 Reward Hacking in Reinforcement Learning | Lil'Log
<https://lilianweng.github.io/posts/2024-11-28-reward-hacking/>