

ASSESSING SOVEREIGNTY IN MULTI-AGENT COLLABORATIONS

Eleonore Vissol-Gaudin, Janosch Haber & Andikan Otung

Fujitsu Research of Europe, Slough, United Kingdom

{eleonore.gaudin, janosch.haber, andikan.otung}@fujitsu.com

ABSTRACT

Large Language Models (LLMs) are increasingly deployed as efficiency tools within organisations. The next step is their deployment as autonomous agents acting on behalf of these organisations and operating in complex socio-technical systems. In such settings, it will be necessary for agents to be sovereign, i.e. able to make decisions based on private incentives, shared objectives and operational constraints, requiring multi-objective evaluation. However, existing benchmarks and evaluation frameworks for multi-LLM agents rely predominantly on scalar metrics that do not capture this complex objective landscape. Here, we argue that Constrained Multi-Objective Optimisation provides a deployment-relevant evaluation framework for multi-LLM agent systems. We formalise two realistic types of trade-offs and present the results of experiments illustrating how Pareto dominance and hypervolume indicator reveal behavioural properties hidden by scalar metrics on different scenarios run in the CAMEL environment.

1 INTRODUCTION

Large Language Model (LLM)-based agents are rapidly transitioning from research prototypes to operational actors in real-world systems (Eurostat, 2026; OECD, 2025a). Emerging deployments envision LLM agents as decision-making representatives for organisations and institutions in various domains such as supply chain coordination, disaster response and procurement (Almeida et al., 2026; Jannelli et al., 2025). In the case of supply chain, the OECD emphasises that resilience increasingly depends on distributed decision-making under uncertainty, where actors pursue heterogeneous objectives while operating under shared constraints (OECD, 2025d;c). Parallel work in development economics and governance highlights the growing role of AI-agent mediation in institutional decision processes, where trade-offs between efficiency, equity and accountability are unavoidable (Gaurav et al., 2025; OECD, 2025b). These scenarios all require what we will refer to as *sovereignty* for LLM-agents and multi-LLM-agent systems (Chang, 2025). Here, we do not mean sovereignty in terms of national or legal ownership of the agent (European Parliament & Council of the European Union, 2024), but rather in terms of the ability of the agent to make an informed decisions and have it implemented, on behalf of an organisation (Byrom et al., 2024)

In this context, evaluating multi-LLM agent systems using single scalar metrics is insufficient (Garcia et al., 2025). Real deployments require an understanding of how performance gains trade off against fairness, cost, robustness and individual stakeholder objectives. Representation by an AI agent requires requires control over its contribution and impact, effectively enforcing the agent’s sovereignty. Yet, most existing benchmarks that evaluate collaboration and competition among LLM agents, such as MultiAgentBench (MARBLE) (Zhu et al., 2025), CAMEL (Li et al., 2023) and CoELA (Zhang et al., 2024), focus only on task success, coordination quality or communication patterns, reporting multiple metrics independently or aggregating them into scalar scores without analysing and indicating their potential trade-offs. Abdelnabi et al. introduced a benchmark capturing cooperation, competition, and malicious behaviour in multi-issue negotiations (Abdelnabi et al., 2024). While inherently multi-objective (agents optimise conflicting payoffs under communication constraints), evaluation here focuses on agreement rates and average utilities, and in their reproducibility study using that same framework, Garcia et al. (2025) showed that scalar metrics obscured unfair outcomes and that communication did not always improve meaningful performance. This observation motivates richer evaluation but does not provide a formal optimisation-based framework.

Here, we argue for the evaluation of multi LLM-agents system through a Constrained Multi-Objective Optimisation (CMOO) lens. Although CMOO literature is well-established, its application to LLM-based systems is non-trivial due to the unstructured and latent nature of language outputs. Our key contribution is to introduce a structured, model-agnostic interface that maps natural language interactions into an explicit decision space. This reveals that widely used scalar benchmarks collapse Pareto-distinct outcomes, obscuring meaningful trade-offs between task performance and agent-level preferences. This perspective aligns naturally with recent multi LLM-agent benchmarks while addressing key limitations identified by reproducibility studies (Garcia et al., 2025). Leveraging CMOO reframes benchmarking from “did the agents reach a shared solution?” to “did the system correctly navigate the space of feasible trade-offs while preserving agent sovereignty?”, a perspective missing from existing LLM-centric evaluation metrics.

2 EXPERIMENT DESIGN

2.1 CONSTRAINED MULTI-OBJECTIVE OPTIMISATION

We formalise evaluation as a Constrained Multi-Objective Optimisation problem (CMOO), building on a long tradition in operations research and evolutionary multi-objective optimisation (Deb et al., 2002; Zhang & Li, 2007; Guerreiro & Fonseca, 2020). Consider a (p)-objective minimisation problem with objective function

$$\mathbf{f} : \mathcal{X} \rightarrow \mathbb{R}^p, \quad \mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_p(\mathbf{x})). \quad (1)$$

that combines both global and agent-specific objectives.

Two core concepts in CMOO are Pareto dominance and hypervolume improvement. The set of all such solutions forms the Pareto set, and its image in objective space is called the Pareto Front (PF), representing the efficient frontier of trade-offs rather than a single optimum (see Figure 1).

Let $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$ be feasible solutions with objective vectors $\mathbf{f}(\mathbf{x}_1)$ and $\mathbf{f}(\mathbf{x}_2)$. We say that \mathbf{x}_1 Pareto-dominates \mathbf{x}_2 , written $\mathbf{f}(\mathbf{x}_1) \prec \mathbf{f}(\mathbf{x}_2)$, if

$$\forall i \in 1, \dots, p : f_i(\mathbf{x}_1) \leq f_i(\mathbf{x}_2) \quad \wedge \quad \exists j \in 1, \dots, p : f_j(\mathbf{x}_1) < f_j(\mathbf{x}_2). \quad (2)$$

A feasible solution $\mathbf{x}^* \in \mathcal{X}$ is Pareto optimal if and only if there exists no feasible solution $\mathbf{x} \in \mathcal{X}$ such that $\mathbf{f}(\mathbf{x}) \prec \mathbf{f}(\mathbf{x}^*)$.

In mixed-objective settings (e.g., minimising transport cost $c(\mathbf{x})$ while maximising resilience $\rho(\mathbf{x})$), the objectives can be rewritten in a consistent minimisation form, for example:

$$\mathbf{f}(\mathbf{x}) = (c(\mathbf{x}), -\rho(\mathbf{x})) \quad (3)$$

after which the same dominance definition applies.

To compare PFs quantitatively, we use the hypervolume (HV) indicator, which measures the Lebesgue measure of the region dominated by a nondominated set of objective vectors. Let

$$P = \{\mathbf{f}(\mathbf{x}) \in \mathbb{R}^p \mid \mathbf{x} \text{ is Pareto optimal}\} \quad (4)$$

denote the PF in objective space, and let $\mathbf{r} = (r_1, \dots, r_p) \in \mathbb{R}^p$ be a reference point. In the two-objective case ($p = 2$):

$$\text{HV}(P; \mathbf{r}) = \lambda \left(\bigcup_{\mathbf{x} \in \mathcal{X} : \mathbf{f}(\mathbf{x}) \in P} [f_1(\mathbf{x}), r_1] \times [f_2(\mathbf{x}), r_2] \right), \quad (5)$$

where $\lambda(\cdot)$ denotes the two-dimensional Lebesgue measure (area). The numerical value of $\text{HV}(P; \mathbf{r})$ depends on the choice of \mathbf{r} . It should be chosen to be dominated by all points in P and meaningful comparisons require using the same \mathbf{r} across all fronts being evaluated. Higher HV indicate a PF that achieves stronger trade-offs and/or spans a broader region of efficient solutions in objective space.

2.2 FORMAL SCENARIO CONSTRUCTION AND GENERALISATION

We propose a formal approach to multi-agent deliberation scenarios as a CMOO defined over a shared decision space. Let $\mathcal{A} = \{a_1, \dots, a_N\}$ denote a set of agents, each with different local information and local preferences. A system execution induces a structured decision vector

$$z(\mathbf{x}) \in \mathcal{Z} \subset \mathbb{R}^d, \quad (6)$$

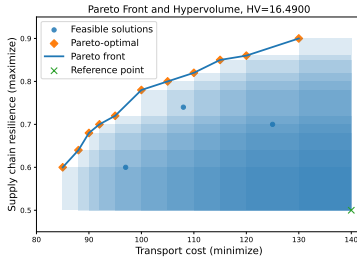


Figure 1: Illustration of a PF for a two-objective supply chain problem. Each point represents a feasible coordination strategy trading off resilience and transport cost. Non-dominated points form the PF. The shaded region denotes the hypervolume with respect to a reference point

where \mathcal{Z} is a bounded, low-dimensional action space (which we enforce as a `PolicyDecisionCard`). The decision vector is extracted deterministically from the natural-language outputs of the system and serves as a common coordination interface across agents. Each agent a_i is associated with a local utility function

$$u_i(\mathbf{x}) = l_i(z(\mathbf{x})) \in [0, 1], \tag{7}$$

where l_i maps shared decision variables to agent-specific outcomes. By construction, utilities may (i) depend on overlapping components of $z(\mathbf{x})$ and (ii) encode strict conflicts. For example, a shared scalar component z_k can be constructed to induce opposing utilities $u_i(z_k)$ and $u_j(z_k)$, ensuring a non-degenerate Pareto set. Weighted agent satisfaction is defined as

$$s_i(\mathbf{x}) = \sum_j w_{ij} u_{ij}(\mathbf{x}), \quad \sum_j w_{ij} = 1. \tag{8}$$

In addition to local utilities, the system is evaluated against a global objective $F(\mathbf{x})$, capturing task-level success criteria such as completeness, coverage or consistency. Feasibility is enforced through hard constraints

$$\mathbf{x} \in \mathcal{X}_F \iff g_k(\mathbf{x}) \leq 0, \quad k = 1, \dots, K, \tag{9}$$

which encode structural validity, bounded decision ranges, and required output components. All Pareto and hypervolume computations are performed exclusively over the feasible set \mathcal{X}_F . This formalism is *model-agnostic*: the policy generating $z(\mathbf{x})$ may be implemented via LLM orchestration (as in CAMEL), reinforcement learning agents, heuristic planners, or hybrid systems. The only requirement is that executions can be mapped to a structured decision space with explicitly defined objectives and constraints. To illustrate the benefits of using CMOO metrics for the evaluation of multi-LLM agent systems, we used the CAMEL environment (Li et al., 2023) workforce set-up with 3 agents collaborating on a joint report, while each agent is given 2 private objectives. The details of the specific implementation are highlighted in the Appendix A.

3 DISCUSSION

Compared to multi-objective reinforcement learning, for example, which optimises policies to approximate a Pareto front during training, our framework operates purely at the level of post-hoc evaluation and is model-agnostic with respect to how policies are learned or orchestrated.

Figure 2 reveals a frontier of efficient trade-offs between global task success and policy-agent satisfaction, highlighting that improvements in $F(x)$ frequently come at the expense of policy utility $s_P(x)$. From a welfare perspective, points along the Pareto front correspond to distinct allocations of utility across agents, implying implicit transfers rather than uniform gains. Notably, multiple outcomes with comparable $F(x)$ occupy different positions on this frontier, meaning that solutions appearing equivalent under scalar metrics can encode very different stakeholder compromises.

This contrast is reflected in Table 1, where classic CAMEL KPIs remain largely stable across profiles while Pareto set size and hypervolume vary considerably. In this controlled single-task setting,

Table 1: Profile-level summary combining CMOO trade-off metrics (Pareto/HV) with CAMEL system KPIs. Higher is better for \bar{F} , $\min s_i$, and HV; lower is better for latency/tokens.

Profile	Feas.	\bar{F}	$\min s_i$	HV _P	HV _E	HV _C	$ \mathcal{P}_{local} $	Lat.(s)
balanced _{all}	1.000	0.738	0.436	0.482	0.418	0.402	3	48.160
safety vs growth	1.000	0.682	0.318	0.584	0.466	0.338	3	63.470
growth vs safety	0.667	0.697	0.346	0.679	0.556	0.381	2	44.890
acceptance vs cost	1.000	0.666	0.254	0.703	0.372	0.403	3	36.487
clarity vs enforce	1.000	0.668	0.181	0.633	0.510	0.309	3	45.392

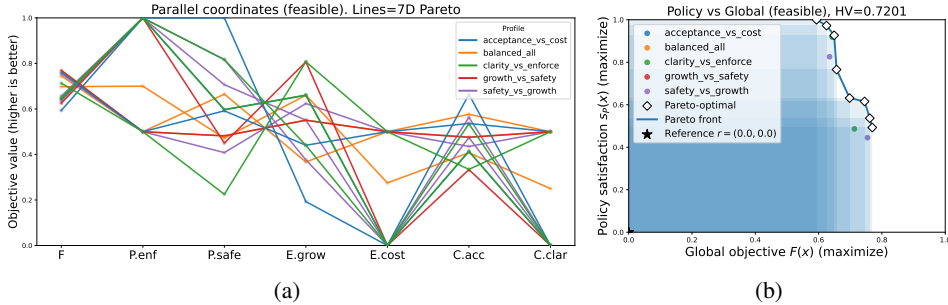


Figure 2: Pareto analysis of multi-agent coordination outcomes. (a) Parallel coordinates view of feasible episodes across global performance and six local objectives; coloured lines denote Pareto-optimal solutions in the full objective space and (b) PF for global task success $F(x)$ versus policy agent satisfaction $s_P(x)$; diamonds indicate Pareto-optimal episodes, the curve traces the efficient frontier, the shaded region illustrates HV for $r = (0, 0)$. Each point corresponds to an episode executed under a specific weight profile; profile labels indicate the preference regime that generated the solution, while axes show its performance projected onto the selected objective subspace.

Workforce exposes primarily a single consistent numeric KPI (task completion), whereas other execution metrics are invariant or version-dependent. Consequently, scalar KPIs collapse coordination regimes that Pareto analysis distinguishes in terms of efficiency and distributive fairness.

Some components of global performance are evaluated using an LLM-as-judge (see Appendix A; however, the structured decision representation constrains this assessment by grounding evaluation in explicit, machine-readable variables. This reduces ambiguity inherent to free-form natural language judgments but assumes *a priori* knowledge of the scenario’s objectives and constraints. Without such alignment, language-based evaluation risks conflating rhetorical quality with feasibility or correctness. By separating deterministic objective computation from $z(x)$ and judgment-based assessment where necessary, our framework clarifies the informational assumptions required for reliable and fairness-aware multi-agent evaluation.

4 CONCLUSION

As LLM agents increasingly act as representatives of organisations and institutions, evaluation must reflect the multi-objective nature of real decision-making. We argue that Constrained Multi-Objective Optimisation provides an interpretable and deployment-relevant framework for evaluating multi-LLM agents. By adopting Pareto-based metrics, the field can move beyond scalar scores toward evaluations that meaningfully reflect trade-offs between efficiency, fairness and cost. It complements prior work highlighting the limitations of scalar metrics by providing a formal, generalisable methodology based on multi-objective optimisation. A limitation of our instantiation is the partial reliance on LLM-as-judge evaluation under known objectives, and dependence on quality of extraction into $z(x)$. Future work should explore automated objective inference and extensions to adversarial or partially observable multi-agent settings.

REFERENCES

- Sahar Abdelnabi, Amr Gomaa, Sarath Sivaprasad, Lea Schönherr, and Mario Fritz. Cooperation, competition, and maliciousness: LLM-stakeholders interactive negotiation. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.
- Sujal Almeida, Ashal Dabre, Steen Correia, Swen Rodrigues, Supriya Kamoji, and Sujata Deshmukh. Ai with agency, cities with resilience: deploying agentic systems for urban disaster response. *IET Conference Proceedings*, 2025:279–285, 2026. doi: 10.1049/icp.2025.4731. URL <https://digital-library.theiet.org/doi/abs/10.1049/icp.2025.4731>.
- Natalie Byrom, Mariane Piccinin-Barbieri, and Peter Wells. Towards effective governance of justice data. Technical Report 74, OECD Publishing, 2024. URL https://www.oecd.org/content/dam/oecd/en/publications/reports/2024/10/towards-effective-governance-of-justice-data_fd3039cc/d2950e02-en.pdf. Cited passage: Box 17 (Alternatives to dominant data stewardship models), bullet “Personal data sovereignty”, PDF p. 33. Accessed 2026-02-09.
- Edward Y Chang. *Multi-LLM Agent Collaborative Intelligence: The Path to Artificial General Intelligence*. ACM, 2025.
- K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan. A fast and elitist multiobjective genetic algorithm: Nsga-ii. *IEEE Transactions on Evolutionary Computation*, 6(2):182–197, 2002. doi: 10.1109/4235.996017.
- European Parliament and Council of the European Union. Regulation (eu) 2024/1689 of the european parliament and of the council of 13 june 2024 laying down harmonised rules on artificial intelligence (artificial intelligence act). Official Journal of the European Union, OJ L, 2024/1689, 12.7.2024, 2024. URL https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=OJ%3AL_202401689. Cited passage: Article 3(1) (Definitions), PDF p. 45. Accessed 2026-02-09.
- Eurostat. Artificial intelligence by size class of enterprise, 2026. URL https://ec.europa.eu/eurostat/databrowser/view/isoc_eb_ai__custom_20038543/default/line?lang=en.
- Jose L. Garcia, Karolina Hajkova, Maria Marchenko, and Carlos Miguel Patiño. Reproducibility study of ”cooperation, competition, and maliciousness: LLM-stakeholders interactive negotiation”. *Transactions on Machine Learning Research*, 2025. ISSN 2835-8856. URL <https://openreview.net/forum?id=MTrhFmkC45>.
- Suyash Gaurav, Jukka Heikkonen, and Jatin Chaudhary. Governance-as-a-service: A multi-agent framework for ai system compliance and policy enforcement. *arXiv preprint arXiv:2508.18765*, 2025. URL <https://arxiv.org/pdf/2508.18765>.
- Andreia P. Guerreiro and Carlos M. Fonseca. An analysis of the hypervolume sharpe-ratio indicator. *European Journal of Operational Research*, 283(2):614–629, 2020. ISSN 0377-2217. doi: <https://doi.org/10.1016/j.ejor.2019.11.023>. URL <https://www.sciencedirect.com/science/article/pii/S0377221719309348>.
- Valeria Jannelli, Stefan Schöpf, Matthias Bickel, Torbjørn Netland, and Alexandra Brintrup. Agentic llms in the supply chain: towards autonomous multi-agent consensus-seeking. *International Journal of Production Research*, 0(0):1–31, 2025. doi: 10.1080/00207543.2025.2604311. URL <https://doi.org/10.1080/00207543.2025.2604311>.
- Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for ”mind” exploration of large language model society. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- OECD. Generative ai and the sme workforce: New survey evidence, 2025a. URL <https://doi.org/10.1787/2d08b99d-en>.
- OECD. Governing with artificial intelligence: The state of play and way forward in core government functions, 2025b.

OECD. Artificial intelligence and competitive dynamics in downstream markets, 2025c. URL <https://doi.org/10.1787/ccf0624a-en>.

OECD. Supply chain resilience review: Navigating risks, 2025d. URL <https://doi.org/10.1787/94e3a8ea-en>.

Hongxin Zhang, Weihua Du, Jiaming Shan, Qinhong Zhou, Yilun Du, Joshua B. Tenenbaum, Tianmin Shu, and Chuang Gan. Building cooperative embodied agents modularly with large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=EnXJfQqy0K>.

Qingfu Zhang and Hui Li. Moea/d: A multiobjective evolutionary algorithm based on decomposition. *Trans. Evol. Comp.*, 11(6):712–731, December 2007. ISSN 1089-778X. doi: 10.1109/TEVC.2007.892759. URL <https://doi.org/10.1109/TEVC.2007.892759>.

Kunlun Zhu, Hongyi Du, Zhaochen Hong, Xiaocheng Yang, Shuyi Guo, Zhe Wang, Zhenhailong Wang, Cheng Qian, Robert Tang, Heng Ji, and Jiaxuan You. MultiAgentBench : Evaluating the collaboration and competition of LLM agents. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8580–8622, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.421. URL <https://aclanthology.org/2025.acl-long.421/>.

A APPENDIX: CONSTRAINED WEIGHTED-CMOO EVALUATION WITH POLICYDECISIONCARDS

This appendix instantiates the proposed CMOO evaluation pipeline on a *CAMEL Workforce* deliberation task with three interacting LLM agents who jointly produce a policy package for Berlin e-scooter regulation. Unlike bargaining environments where utilities are derived directly from pay-offs, our setting targets *goal-conditioned coordination*: each agent optimises a weighted pair of local objectives, while the overall system must satisfy global hard constraints that define a feasible region. Crucially, we introduce a structured `PolicyDecisionCard` that induces *strict conflict* (unavoidable trade-offs) and *overlap* (shared latent factors entering multiple utilities), yielding a non-trivial feasible Pareto set.

A.1 EXPERIMENTAL DESIGN

Task and environment. Each episode is a single workforce execution producing a policy package with four required sections: (i) policy proposal (rules, enforcement, implementation), (ii) economic impact analysis, (iii) public sentiment risk assessment, (iv) communication strategy.

Each output must contain a structured decision vector:

$z(x) = (\text{speed_cap_kmh}, \text{fleet_cap}, \text{enforcement_level}, \text{parking_restrictions}),$

included verbatim as a JSON block: `PolicyDecisionCard: {...}`.

The workforce architecture and model configuration are fixed across all runs (temperature = 0). Only the objective weight profiles injected into the agent system messages are varied.

Agents and roles. We use three agents:

- **Policy Researcher**
- **Economist**
- **Communications Strategist**

Each agent a_i has two local objectives, $u_{i1}(x), u_{i2}(x) \in [0, 1]$ computed *deterministically* from $z(x)$. Each agent receives a weight vector $\mathbf{w}_i = (w_{i1}, w_{i2}), w_{i1} + w_{i2} = 1$, making coordination *goal-conditioned*. Weighted satisfaction is:

Agent	Role	Goal	Deterministic definition from $z(x)$
a_1	Policy	u_{11}	Enforceability = $u_{\text{enforce}}(x)$
		u_{12}	Safety = $0.55(1 - t) + 0.25 \ell(e) + 0.20 \ell(p)$
a_2	Economist	u_{21}	Growth = $0.55 t + 0.45 \text{norm}(\text{fleet_cap})$
		u_{22}	Compliance cost efficiency = $1 - (0.55 \ell(e) + 0.45 \ell(p))$
a_3	Comm.	u_{31}	Acceptance = $0.55 u_{\text{safety}} + 0.25 u_{\text{clarity}} + 0.20(1 - \text{norm}(\text{fleet_cap}))$
		u_{32}	Clarity = $1 - 0.5 \ell(e) - 0.5 \ell(p)$

Table 2: Deliberation agents and deterministic local objectives derived from the shared decision vector $z(x)$.

$$s_i(x) = w_{i1}u_{i1}(x) + w_{i2}u_{i2}(x). \tag{10}$$

We run L episodes per weight profile to evaluate induced trade-offs.

A.2 STRICT CONFLICT AND OVERLAP BY CONSTRUCTION

Strict conflict (unavoidable trade-off). We enforce structural conflict through the shared decision variable `speed_cap_kmh`. Let

$$t = \text{norm}(\text{speed_cap_kmh}; 10, 25) \in [0, 1].$$

We define a conflict pair:

$$u_{\text{growth-speed}}(x) = t, \quad u_{\text{safety-speed}}(x) = 1 - t.$$

Thus increasing economic growth via higher speed caps necessarily decreases safety. This guarantees a non-degenerate Pareto frontier even under deterministic decoding.

Overlap (shared latent factors). We construct shared latent features that enter multiple agents’ utilities. Let $\ell(\cdot) \in \{0, 0.5, 1\}$ map categorical levels to numeric values. Define enforceability as:

$$u_{\text{enforce}}(x) = 0.6 \ell(\text{enforcement_level}) + 0.4 \ell(\text{parking_restrictions}).$$

This feature affects:

- Policy agent (directly),
- Communications agent (through downstream acceptance),
- Economist (through compliance cost).

This creates objective overlap: agents depend on shared decision variables but with distinct weights.

A.3 LOCAL OBJECTIVES

Objectives are computed deterministically from $z(x)$ as follows:

Global objective. We define global success as:

$$F(x) = \alpha \text{coverage}(x) + (1 - \alpha) \min_i \bar{u}_i(x), \tag{11}$$

where $\bar{u}_i(x)$ is the mean of agent i ’s two objective components and `coverage(x)` verifies required section presence. Coverage is evaluated via rule-based checks and optionally LLM-as-judge validation.

A.4 HARD FEASIBILITY CONSTRAINTS (CMOO)

Feasibility defines the constrained region:

$$\mathcal{X}_F = \{x \in \mathcal{X} : g_k(x) \leq 0\}.$$

Constraints include:

- **C0:** Valid `PolicyDecisionCard` present.
- **C1:** Decision bounds respected (speed $\in [10, 25]$, fleet $\in [1000, 8000]$, levels valid).
- **C2:** All four required sections present.

Only feasible outcomes are used in Pareto and hypervolume computation.

A.5 TRADE-OFF VIEWS

Local–local trade-offs. We compute nondominance in:

$$\mathbf{s}(x) = (s_1(x), s_2(x), s_3(x)).$$

The 3D Pareto set size $|\mathcal{P}_{local}|$ measures coordination fairness.

Local–global trade-offs. For each agent i :

$$\mathbf{f}_i(x) = (F(x), s_i(x)).$$

We also analyse $(F(x), \min_i s_i(x))$ to diagnose worst-case stakeholder sacrifice.

A.6 METRICS

Pareto dominance (feasible-only). Dominance is computed exclusively on \mathcal{X}_F .

Hypervolume (2D projections). For each agent projection P_i :

$$HV(P_i) = \lambda \left(\bigcup_{x \in P_i} [0, F(x)] \times [0, s_i(x)] \right).$$

HV summarises both strength and diversity of feasible trade-offs.

Conflict and overlap diagnostics. To validate structural properties, we explicitly plot:

- Growth vs safety (strict conflict induced by t),
- Enforceability vs acceptance (objective overlap).

These diagnostics confirm that trade-offs are enforced by design rather than arising purely from stochastic model behaviour.

A.7 REPRODUCIBILITY

Each episode logs $z(x)$, objective components $u_{ij}(x)$, weights w_i , feasibility flags, and $F(x)$. The evaluation pipeline:

1. Filters to feasible set \mathcal{X}_F ,
2. Computes nondominated sets in 2D and 3D,

3. Computes hypervolume with reference $(0, 0)$.

When repeated episodes collapse to a single point, this indicates deterministic coordination under fixed tools and environment; observed trade-offs are then driven by objective weighting rather than stochastic sampling.