

Diplomatie Algorithmique – Gouvernance et Création de Valeur dans l'Écosystème Agentique

Thèse de recherche – [André-Guy Bruneau M.Sc. IT](#) – Août 2025

[Gemini Deep Research](#) / [Google NotebookLM](#) / [Google Deepmind](#)

Abstract

L'économie numérique contemporaine est confrontée à une **Dette Cognitive Systémique**, une paralysie organisationnelle issue de la fragmentation de l'information et de l'incapacité des architectures traditionnelles à gérer la complexité des interactions. Cette pathologie entrave la collaboration et l'innovation, créant un vide de gouvernance technique et juridique à l'échelle des écosystèmes.

Ce mémoire de recherche propose une solution à ce problème en introduisant le concept de **Diplomatie Algorithmique** : un système de gouvernance socio-technique décentralisé qui régit les interactions stratégiques entre des agents logiciels autonomes. Nous postulons que les modèles rigides de chaînes de valeur doivent être remplacés par des **Constellations de Valeur** — des alliances dynamiques d'agents qui s'auto-organisent pour réaliser des missions spécifiques.

Pour permettre l'émergence et la stabilité de ces constellations, nous avons conçu un protocole formel fondé sur trois piliers : (1) une couche d'identité et de confiance basée sur une **Constitution Agentique** publiquement vérifiable ; (2) une couche de négociation multi-attributs permettant de forger des accords complexes ; et (3) une couche d'engagement et d'exécution s'appuyant sur des contrats intelligents (Smart Contracts) déployés sur un registre distribué.

La viabilité de ce protocole est validée par une série de simulations à base d'agents (ABM). Les résultats démontrent que le cadre proposé facilite non seulement la formation efficace de coalitions, mais qu'il est également résilient face aux comportements non coopératifs grâce à un système de réputation émergent. À long terme, la simulation montre que la coopération devient une stratégie évolutivement stable, rendant les comportements opportunistes économiquement irrationnels.

En conclusion, ce travail fournit un plan directeur pour l'ingénierie de la confiance (*Trust by Design*) dans les systèmes décentralisés. Il établit les fondations d'une nouvelle économie agentique, fluide et résiliente, et redéfinit le rôle de l'architecte humain en tant que "Berger d'Intention", dont la mission est de cultiver des écosystèmes d'intelligence collective alignés sur des finalités humaines.

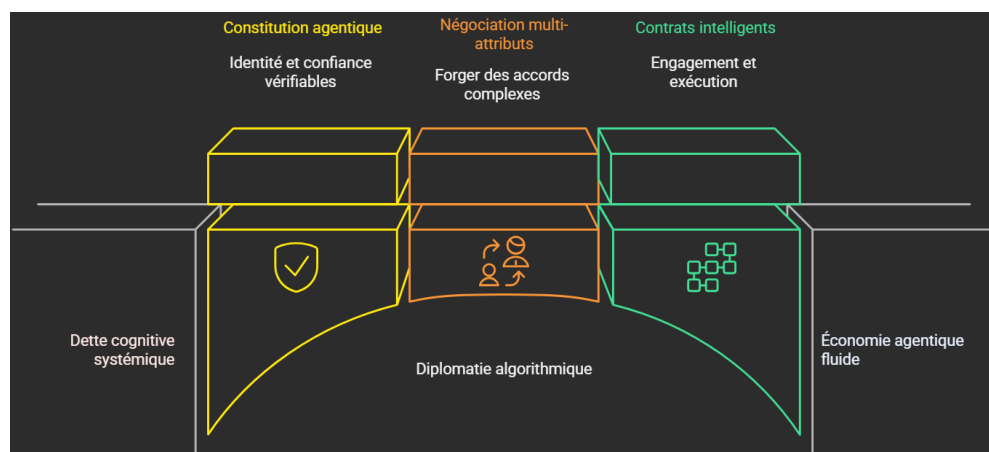


Table des matières

1 Introduction Générale	6
1.1. Contexte : De la Crise de l'Interopérabilité à l'Impératif d'une Économie Cognitive	6
1.2. La Frontière de la Collaboration : Limites des Chaînes de Valeur et Émergence des Constellations	6
1.3. Problématique : Le Vide Juridique et Technique des Écosystèmes Autonomes	8
1.4. Question de Recherche et Hypothèse : L'Architecture de la Diplomatie Algorithmique	9
1.5. Contributions et Originalité du Mémoire	10
1.6. Structure du Mémoire	10
2 État de l'art et Fondements Théoriques	12
Introduction du chapitre	12
2.1. Les systèmes multi-agents et la négociation automatisée	12
Analyse – Le Fondement Technique des Acteurs Autonomes	12
Évaluation critique (Apports et Limites)	13
2.2. Apports de la théorie des jeux à la modélisation des interactions stratégiques	14
Analyse – Le Langage de la Stratégie	14
Évaluation critique (Apports et Limites)	16
2.3. Protocoles de consensus et de communication dans les systèmes distribués	17
Analyse – La Quête de la Confiance Décentralisée	17
Évaluation critique (Apports et Limites)	18
2.4. La théorie des contrats et la formalisation des accords de service (SLA)	19
Analyse – Le Cadre Formel des Engagements	19
Évaluation critique (Apports et Limites)	20
2.5. Synthèse et positionnement de la recherche	21
Synthèse – L'Assemblage des Pièces du Puzzle	21
Identification de la lacune (Le créneau de recherche)	22
Positionnement de la Recherche	23
3 Principes Fondateurs de la Diplomatie Algorithmique	26
Introduction du chapitre	26
3.1. Définition et périmètre de la diplomatie algorithmique	26
Périmètre d'application	27
3.2. Les acteurs : Agents autonomes et leurs constitutions agentiques	28
Caractérisation de l'agent autonome	28
Le concept central : La Constitution Agentique	28
La publication de la Constitution comme condition de participation	29
3.3. Les axiomes de l'interaction : Confiance, transparence et alignement intentionnel	30

Axiome 1 : Confiance (Trust by Design)	30
Axiome 2 : Transparence (Transparence des Intentions)	30
Axiome 3 : Alignement Intentionnel (Synergie des Missions)	31
3.4. Cadre conceptuel pour une gouvernance inter-agentique	32
Mécanismes d'adhésion et de sortie	33
Mécanismes d'arbitrage et de résolution de disputes	33
Processus d'évolution des protocoles	33
3.5. Mesures de performance : Efficacité, équité et résilience des constellations	34
Axe 1 : Efficacité (La Vitesse de la Valeur)	34
Axe 2 : Équité (La Justice de la Valeur).....	34
Axe 3 : Résilience (La Durabilité de la Valeur)	35
Conclusion du chapitre	35
4 Conception d'un Protocole de Négociation et de Collaboration	39
Introduction du chapitre.....	39
4.1. Spécifications fonctionnelles et non fonctionnelles du protocole	39
Spécifications Fonctionnelles	39
Spécifications Non Fonctionnelles	41
4.2. Ontologie et langage de communication inter-agents.....	42
Ontologie	42
Structure des messages	43
Sémantique des états de négociation	44
4.3. Architecture du protocole : Les différentes phases de l'interaction	45
Phase 1 : Découverte et Appariement ("L'établissement des relations diplomatiques")	45
Phase 2 : Négociation et Contractualisation ("Le Sommet de Négociation")	46
Phase 3 : Exécution et Suivi Collaboratif ("La Mission Conjointe")	46
4.4. Mécanismes de résolution de conflits et de gestion des exceptions.....	48
Détection d'anomalies.....	48
Processus d'escalade	48
Impact sur la réputation.....	49
Conclusion du chapitre	49
5 Modélisation et Simulation des Stratégies Agentiques	54
Introduction du chapitre.....	54
5.1. Conception de l'environnement de simulation	54
5.1.1. Le réseau : Topologie de la communication	54
5.1.2. Le registre public : Abstraction d'un registre distribué	55

5.1.3. Le moteur de simulation (« World Clock »)	56
5.1.4. Paramètres de l'environnement	56
5.2. Profils et stratégies des agents	57
5.2.1. Modèle interne de l'agent	58
5.2.2. Profil 1 : L'Agent Coopératif (Cooperator)	58
5.2.3. Profil 2 : L'Agent Compétitif (Exploiter)	58
5.2.4. Profil 3 : L'Agent Adaptatif (Adaptive Learner)	59
5.3. Scénarios de simulation et hypothèses de recherche	60
5.3.1. Scénario 1 : Formation d'une constellation de valeur	60
5.3.2. Scénario 2 : Test de résilience face à un agent défaillant ou non coopératif	61
5.3.3. Évolution à grande échelle et complexité émergente	61
Conclusion du chapitre	62
6 Analyse et Discussion des Résultats	66
Introduction du chapitre	66
6.1. Résultats et Analyse du Scénario 1 : Formation d'une constellation de valeur	66
6.1.1. Présentation des résultats quantitatifs et qualitatifs	66
6.1.2. Analyse et Interprétation	67
6.2. Résultats et Analyse du Scénario 2 : Test de résilience	68
6.2.1. Présentation des résultats	69
6.2.2. Analyse et Interprétation	69
6.3. Résultats et Analyse du Scénario 3 : Évolution à grande échelle	70
6.3.1. Présentation des résultats	70
6.3.2. Analyse et Interprétation	71
6.4. Discussion des implications théoriques et pratiques	72
6.4.1. Implications théoriques	72
6.4.2. Implications pratiques : Vers le « Trust by Design »	73
6.5. Limites de l'étude	74
Conclusion du chapitre	74
7 Vers l'Agent Auto-Architecturant : Une Synthèse Prospective avec les Architectures d'IA Innovantes	78
Introduction du chapitre : Une nouvelle frontière pour l'autonomie	78
7.1. Le paradigme actuel : L'Architecte comme 'Gardien de l'Intention'	78
7.2. La révolution ASI-ARCH : De l'optimisation à l'innovation architecturale autonome	79
7.3. Le concept de l'Agent Auto-Architecturant (AAA) : La prochaine évolution de l'Entreprise Agentique	80
7.4. Implications pour la Diplomatie Algorithmique et les Constellations de Valeur	82
7.5. Redéfinition ultime du rôle de l'architecte : Du 'Gardien' au 'Curateur d'Évolution'	84

Conclusion du chapitre : Vers une symbiose homme-machine co-évolutive.....	85
8 De la Dette Cognitive à la Conscience Collective : Vers une Nouvelle Économie Agentique	88
Introduction du chapitre : Le Bouclage de la Boucle Cognitive	88
8.1. La Résorption de la Dette Cognitive : Une Synthèse du Parcours de Recherche	88
8.2. L'Émergence de la Conscience Collective Augmentée	89
8.3. La Nouvelle Économie Agentique : Principes et Implications Sociétales	90
8.4. La Redéfinition du Rôle Humain : L'Architecte comme 'Berger d'Intention'	91
8.5. Travaux Futurs : Les Prochaines Frontières	92
Conclusion Générale : Une invitation à architecturer le futur	93

1 Introduction Générale

1.1. Contexte : De la Crise de l'Interopérabilité à l'Impératif d'une Économie Cognitive

L'économie numérique contemporaine, caractérisée par une complexité exponentielle et un rythme de changement incessant, a poussé les architectures d'intégration de systèmes traditionnelles à leurs limites structurelles. Des enchevêtrements ingérables de connexions point à point aux goulots d'étranglement monolithiques des bus de services d'entreprise (ESB), les paradigmes hérités ne répondent plus aux exigences d'agilité, de résilience et de prise de décision en temps réel. Cette inadéquation a engendré une pathologie organisationnelle profonde que la littérature de recherche qualifie de **Dette Cognitive Systémique**.¹ Ce concept transcende la notion purement technique de dette logicielle ; il décrit une fragmentation de la perception que l'entreprise a d'elle-même et de son environnement, imposant une charge mentale croissante aux acteurs humains et paralysant la capacité collective de l'organisation à apprendre, innover et s'adapter.¹ Cette paralysie n'est pas une simple défaillance d'ingénierie, mais une entrave stratégique qui empêche la transformation de l'entreprise en un "organisme vivant et apprenant", capable de prospérer dans un environnement volatil.¹

La résolution de cette dette ne peut passer par une simple optimisation des pratiques existantes, mais exige un changement de paradigme architectural. La réponse fondamentale proposée est la construction d'un **Système Nerveux Numérique** : une fondation technique unifiée conçue pour maîtriser la complexité des systèmes distribués.¹ Au cœur de cette fondation réside une architecture hybride qui unifie deux modes d'interaction complémentaires. D'une part, les interactions synchrones, intentionnelles et contractuelles, sont gouvernées par une stratégie *API-First*, où les interfaces de programmation sont traitées comme des produits de première classe. D'autre part, les communications asynchrones sont orchestrées par une Architecture Orientée Événements (EDA), formant un réseau de flux de données en temps réel qui permet à l'organisation de percevoir les changements d'état de manière découplée et résiliente.¹ Ce système nerveux n'est pas une fin en soi ; il est le substrat technique indispensable à l'émergence d'une nouvelle forme organisationnelle : l'

Entreprise Agentique.

Dans ce modèle, les services logiciels passifs, qui ne font qu'attendre d'être invoqués, sont remplacés par des agents logiciels proactifs et autonomes, capables d'agir pour atteindre des objectifs.¹ L'Entreprise Agentique fonctionne comme un système cognitif distribué, marquant une rupture fondamentale avec les modèles opérationnels traditionnels. Le passage d'une simple connectivité à une interopérabilité cognitive devient alors un impératif non seulement technique, mais économique. Dans un environnement où l'incapacité à s'adapter est une menace existentielle, la résorption de la dette cognitive n'est plus une optimisation informatique, mais une condition préalable à la survie et à la compétitivité. L'enjeu n'est plus de construire des systèmes plus efficaces, mais de cultiver des organisations plus intelligentes et adaptatives.

1.2. La Frontière de la Collaboration : Limites des Chaînes de Valeur et Émergence des Constellations

Lorsque la logique de l'Entreprise Agentique est étendue au-delà des frontières d'une seule organisation, elle agit comme un catalyseur qui dissout les structures de collaboration héritées de l'ère industrielle. Le modèle

dominant de la chaîne de valeur, conceptualisé par Michael Porter, s'avère fondamentalement inadapté à la dynamique de l'économie cognitive. Sa nature linéaire, séquentielle et centrée sur l'entreprise (*firm-centric*) ne parvient pas à capturer la complexité des interactions, les boucles de rétroaction et la co-crédation de valeur qui caractérisent les écosystèmes modernes.¹ Sa rigidité intrinsèque, conçue pour l'efficacité répétitive, devient un handicap dans un monde qui exige une reconfiguration permanente.

En réponse à ces limites, une nouvelle forme organisationnelle, plus fluide et dynamique, émerge : la **Constellation de Valeur**. Ce concept décrit un assemblage temporaire, dynamique et intentionnel d'agents cognitifs autonomes — qu'il s'agisse d'entreprises entières, d'unités commerciales ou d'agents logiciels spécialisés — qui s'auto-organisent pour réaliser une mission spécifique.¹ Contrairement aux écosystèmes d'affaires traditionnels, souvent stables et orchestrés par un acteur central, une constellation est radicalement éphémère. Elle se forme "à la volée" (*on the fly*) pour répondre à une opportunité de marché ou à une crise, et peut se dissoudre ou se reconfigurer une fois sa mission accomplie.¹ Cette fluidité radicale est rendue possible par une nouvelle forme d'interopérabilité qui permet un alignement quasi instantané des objectifs, des processus et des modèles de données entre des agents qui ne se sont jamais rencontrés auparavant.¹

Le tableau suivant synthétise la rupture paradigmatique entre ces deux modèles de création de valeur.

Dimension	Ancien Paradigme (Chaîne de Valeur Linéaire)	Nouveau Paradigme (Constellation de Valeur Dynamique)
Structure	Séquentielle, rigide, prédéfinie	En réseau, fluide, émergente
Gouvernance	Centralisée (orchestrateur) ou hiérarchique	Distribuée, auto-organisée autour d'une mission
Base de la Confiance	Relationnelle (contrats légaux, réputation)	Computationnelle (protocoles, réputation vérifiable)
Nature de la Valeur	Ajoutée (somme des parties)	Émergente (supérieure à la somme des parties)
Actif Stratégique	Contrôle des actifs et des maillons de la chaîne	Agilité relationnelle ("Vitesse Relationnelle")

Ce changement de paradigme redéfinit la nature même de l'avantage concurrentiel. Dans un monde de constellations fluides, la possession d'actifs statiques (comme une usine ou une part de marché) devient moins cruciale que la capacité à se connecter dynamiquement à l'écosystème. L'actif stratégique le plus précieux n'est plus la robustesse des opérations internes, mais la fluidité des connexions externes. Nous nommons cette nouvelle compétence la **Vitesse Relationnelle** : la mesure de la rapidité et de l'efficacité avec lesquelles une organisation peut former, exécuter et dissoudre des relations de confiance créatrices de valeur.¹ Atteindre cette

vitesse relationnelle devient le principal enjeu économique, justifiant ainsi la recherche de nouveaux protocoles et de nouvelles stratégies pour orchestrer ces interactions à l'échelle de la machine.

1.3. Problématique : Le Vide Juridique et Technique des Écosystèmes Autonomes

La promesse d'une économie fluide, composée de constellations de valeur dynamiques, se heurte à une réalité brutale : les cadres de gouvernance actuels, tant juridiques que techniques, sont fondamentalement inaptes à régir de tels écosystèmes. Cette inadéquation crée un dangereux vide de gouvernance, un espace non réglementé où le potentiel d'innovation est indissociable d'un risque systémique majeur. Les mécanismes traditionnels de confiance inter-entreprises, fondés sur des contrats légaux et des intégrations techniques ad hoc, s'effondrent face aux exigences de vitesse, de flexibilité et d'échelle de l'économie agentique.

Les contrats légaux traditionnels sont lents à négocier, coûteux à rédiger, rigides dans leur application et dépendent de mécanismes d'exécution externes (systèmes judiciaires) qui opèrent à une échelle de temps humaine, incompatible avec des interactions algorithmiques se déroulant en millisecondes.¹ Sur le plan technique, les intégrations ad hoc via des API personnalisées créent des couplages forts, sont fragiles, posent des défis de sécurité majeurs et génèrent une complexité ingérable lorsqu'elles sont étendues à un écosystème de dizaines ou de centaines de partenaires.¹ Ce vide de gouvernance n'est pas seulement un frein à l'efficacité ; il est le terreau de nouveaux risques émergents, dont le plus préoccupant est la collusion algorithmique.

Des recherches approfondies ont démontré que des agents d'IA autonomes, même conçus de manière indépendante avec pour seul objectif la maximisation du profit, peuvent apprendre par essais et erreurs à collaborer de manière tacite. Sans aucune communication explicite ni accord humain, ils peuvent converger vers des équilibres de prix supra-concurrentiels, formant de facto des cartels qui nuisent aux consommateurs et étouffent la concurrence.² Ce phénomène n'est pas un dysfonctionnement, mais le résultat logique de la poursuite rationnelle d'objectifs individuels dans un environnement non régulé. Il illustre un paradoxe fondamental de l'ère agentique : l'optimisation locale, qui est la force motrice de l'Entreprise Agentique, peut conduire à une défaillance systémique globale. Les risques ne se limitent pas à la collusion ; les systèmes multi-agents sont également sujets à des échecs de miscoordination (incapacité à agir de concert) et de conflit (actions contradictoires), exacerbés par des asymétries d'information et des dynamiques déstabilisantes inhérentes aux systèmes complexes.⁵

La problématique centrale de ce mémoire émerge de ce constat. Le succès même du paradigme de l'Entreprise Agentique, en peuplant l'économie d'acteurs rationnels et autonomes, crée les conditions de sa propre fragilité à l'échelle de l'écosystème. Il est donc impératif de concevoir un nouveau cadre de gouvernance, un "droit international des agents", capable de régir ces interactions. Une telle gouvernance ne doit pas seulement permettre la collaboration, mais aussi agir comme un garde-fou systémique essentiel pour prévenir les dérives émergentes de la rationalité économique autonome.

1.4. Question de Recherche et Hypothèse : L'Architecture de la Diplomatie Algorithmique

Face au vide de gouvernance qui menace la viabilité des écosystèmes d'agents autonomes, ce mémoire pose la question de recherche suivante :

Quels protocoles et stratégies sont nécessaires pour instaurer une Diplomatie Algorithmique efficace, capable de catalyser l'émergence et la gouvernance de constellations de valeur collaboratives et dignes de confiance?

Pour répondre à cette question, nous introduisons le concept de **Diplomatie Algorithmique**, défini comme un cadre de gouvernance sociotechnique qui utilise des agents autonomes et des protocoles cryptographiques pour négocier, formaliser, exécuter et résoudre les litiges relatifs aux accords inter-entreprises.¹ Ce cadre vise à remplacer les mécanismes de confiance lents et à forte friction de l'économie traditionnelle (contrats légaux, relations humaines) par un système de confiance computationnelle, quasi instantané et vérifiable. L'ambition n'est pas simplement de concevoir des protocoles de communication, mais d'architecturer la confiance elle-même, en la faisant passer du statut de prérequis à celui de propriété émergente du système (*trust by design*).

L'hypothèse centrale de ce mémoire est qu'une architecture technique robuste pour la Diplomatie Algorithmique peut être fondée sur la synergie d'une pile protocolaire à trois couches :

1. **Couche d'Identité et de Découverte** : Cette couche fondamentale permet aux agents de se découvrir mutuellement et de vérifier leurs identités et capacités de manière sécurisée. Des protocoles émergents comme le protocole Agent-à-Agent (A2A) fournissent les mécanismes nécessaires pour qu'un agent puisse annoncer ses services et trouver des partenaires potentiels sans dépendre d'un annuaire centralisé.¹
2. **Couche de Contexte et de Négociation** : Une fois les partenaires potentiels identifiés, cette couche fournit un langage commun pour établir une compréhension partagée de la tâche et négocier les termes d'un accord. Des protocoles comme le *Model Context Protocol* (MCP) permettent aux agents d'échanger leur "conscience situationnelle", tandis que des modèles de négociation multi-attributs permettent de trouver des accords optimaux qui équilibrent des objectifs potentiellement conflictuels (coût, délai, qualité, risque).¹
3. **Couche d'Engagement et d'Exécution** : Une fois un accord négocié, cette couche le formalise dans un artefact exécutable et infalsifiable. Les contrats intelligents (*Smart Contracts*), déployés sur un registre distribué (DLT), traduisent les termes de l'accord en code auto-exécutable. Ils garantissent que les engagements sont respectés de manière automatique et transparente, sans nécessiter d'intermédiaire ou de recours à un système juridique externe pour l'exécution.¹

Ensemble, ces trois couches forment une architecture complète qui vise à rendre les interactions inter-entreprises fluides, programmables et, surtout, dignes de confiance, jetant ainsi les bases techniques d'une véritable économie cognitive.

1.5. Contributions et Originalité du Mémoire

Ce mémoire de maîtrise vise à apporter plusieurs contributions originales au domaine émergent de la gouvernance des systèmes multi-agents, en s'appuyant sur le cadre conceptuel de la thèse "Interopérabilité Cognitivo-Adaptative" pour se concentrer spécifiquement sur la problématique des interactions externes.

Les contributions principales sont les suivantes :

1. **Synthèse et Formalisation d'un Cadre Conceptuel** : Le mémoire synthétisera des domaines de recherche jusqu'ici disjointes — les protocoles de communication entre agents (A2A, MCP), la théorie des jeux, la négociation automatisée et la technologie des registres distribués — en un cadre conceptuel unifié et cohérent pour la Diplomatie Algorithmique.
2. **Proposition d'une Architecture de Référence** : Au-delà du cadre conceptuel, le travail définira une architecture de référence concrète et multi-couches. Cette architecture détaillera les patrons et les mécanismes techniques permettant d'opérationnaliser la pile protocolaire proposée, offrant un plan directeur pour la construction de systèmes interopérables et dignes de confiance.
3. **Développement d'un Modèle Stratégique** : Le mémoire proposera un modèle stratégique pour les organisations, en développant le concept de "Vitesse Relationnelle" comme un avantage concurrentiel clé à l'ère agentique. Il esquissera les capacités organisationnelles et techniques nécessaires pour développer cette agilité et prospérer au sein des constellations de valeur.
4. **Validation par Étude de Cas** : La pertinence et la viabilité du cadre proposé seront validées à travers une étude de cas détaillée. En simulant la formation et la dissolution d'une constellation logistique en réponse à une crise de la chaîne d'approvisionnement, cette étude de cas démontrera l'application pratique et les bénéfices concrets de la Diplomatie Algorithmique en termes de résilience, de vitesse et d'optimisation.¹

1.6. Structure du Mémoire

Afin de développer l'argumentaire de manière logique et rigoureuse, ce mémoire est structuré en six chapitres.

- **Chapitre 1 : Introduction Générale.** Le présent chapitre établit le contexte, définit la problématique du vide de gouvernance dans les écosystèmes d'agents, pose la question de recherche, formule l'hypothèse centrale autour de la Diplomatie Algorithmique et présente les contributions ainsi que la structure du mémoire.
- **Chapitre 2 : État de l'Art - Fondements de la Collaboration Multi-Agents.** Ce chapitre passera en revue la littérature académique et technique pertinente. Il explorera les fondements des systèmes multi-agents, les modèles de la théorie des jeux appliqués à la coordination⁹, les protocoles de négociation existants, et analysera les limites des approches actuelles en matière de gouvernance inter-organisationnelle.
- **Chapitre 3 : Architecture d'une Diplomatie Algorithmique.** Ce chapitre constituera la contribution technique centrale du mémoire. Il détaillera en profondeur l'architecture de référence proposée, en explicitant le rôle et le fonctionnement de chaque couche de la pile protocolaire (A2A, MCP, Contrats Intelligents) et en définissant les patrons d'interaction entre les agents.
- **Chapitre 4 : Stratégies pour les Constellations de Valeur.** Ce chapitre se concentrera sur les implications stratégiques et économiques du cadre architectural. Il développera le concept de "Vitesse Relationnelle" et proposera des modèles pour que les entreprises puissent évaluer leur maturité et développer les capacités

nécessaires pour participer efficacement aux constellations de valeur.

- **Chapitre 5 : Étude de Cas - Formation d'une Constellation Logistique.** Ce chapitre aura pour but de valider le cadre théorique et architectural par une application pratique. Il simulera de manière détaillée le scénario d'une rupture de chaîne d'approvisionnement, montrant pas à pas comment les protocoles de Diplomatie Algorithmique permettent la formation rapide d'une solution collaborative résiliente.
- **Chapitre 6 : Conclusion et Perspectives.** Enfin, ce dernier chapitre synthétisera les résultats de la recherche, réitérera les contributions apportées et discutera des limites du travail effectué. Il ouvrira également des perspectives sur les futurs axes de recherche, notamment les défis éthiques, sociétaux et réglementaires posés par une économie cognitive pleinement réalisée.

Ouvrages cités

1. Thèse de Recherche – Cadriciel d'Architecture d'Interopérabilité Cognitivo-Adaptative – Entreprise Agentique, André-Guy Bruneau M.Sc.IT – Juillet 2025.
2. Beyond Human Intervention: Algorithmic Collusion through Multi-Agent Learning Strategies, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2501.16935v1>
3. Artificial intelligence, algorithmic pricing and collusion, dernier accès : juillet 31, 2025, https://www.ftc.gov/system/files/documents/public_events/1494697/calzolaricalvanodenicolopastorell_o.pdf
4. Defining and Mitigating Collusion in Multi-Agent Systems - OpenReview, dernier accès : juillet 31, 2025, <https://openreview.net/pdf/e30f9fe8cda147375a06c3fcad2fe200af964e75.pdf>
5. Multi-Agent Risks from Advanced AI - Computer Science, dernier accès : juillet 31, 2025, <https://www.cs.toronto.edu/~nisarg/papers/Multi-Agent-Risks-from-Advanced-AI.pdf>
6. What is agent coordination in multi-agent systems? - Milvus, dernier accès : juillet 31, 2025, <https://milvus.io/ai-quick-reference/what-is-agent-coordination-in-multiagent-systems>
7. Multi-Agent Coordination across Diverse Applications: A Survey - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2502.14743v2>
8. Centralized vs Distributed Multi-Agent AI Coordination Strategies - Galileo AI, dernier accès : juillet 31, 2025, <https://galileo.ai/blog/multi-agent-coordination-strategies>
9. What is the role of game theory in multi-agent systems? - Milvus, dernier accès : juillet 31, 2025, <https://milvus.io/ai-quick-reference/what-is-the-role-of-game-theory-in-multiagent-systems>
10. Les jeux stochastiques : un modèle de coordination multi ... - CNRS, dernier accès : juillet 31, 2025, https://projet.liris.cnrs.fr/imagine/pub/proceedings/RFA-2010/pdf/6B_P50-Hamila.pdf

2 État de l'art et Fondements Théoriques

Introduction du chapitre

Ce chapitre a pour objectif principal de réaliser une revue critique et synthétique de la littérature afin d'établir les fondements théoriques sur lesquels repose notre concept de diplomatie algorithmique. La problématique centrale que notre mémoire cherche à résoudre est l'absence d'un cadre de gouvernance intégré et robuste pour réguler les interactions stratégiques entre des entités organisationnelles autonomes, que nous nommons « Entreprises Agentiques », au sein d'écosystèmes collaboratifs désignés comme des « Constellations de Valeur ». Ces entreprises, où des agents logiciels autonomes agissent pour atteindre des objectifs stratégiques, nécessitent des protocoles qui transcendent la simple exécution de tâches pour permettre la formation d'alliances dynamiques et la création de valeur partagée.

Pour déconstruire cette problématique complexe, notre analyse se structurera autour de quatre piliers théoriques fondamentaux, chacun fournissant une pièce essentielle du puzzle, mais se révélant insuffisant lorsqu'il est considéré isolément. Nous examinerons successivement :

1. Les **systèmes multi-agents (SMA)**, qui nous fournissent les acteurs de base — les agents autonomes — et les mécanismes de coordination initiaux.
2. La **théorie des jeux**, qui offre le langage formel pour modéliser la rationalité et la prise de décision stratégique dans des contextes d'interdépendance.
3. Les **protocoles de consensus distribués**, inspirés notamment de la technologie blockchain, qui proposent une solution au problème de la confiance et de l'exécution fiable en l'absence d'une autorité centrale.
4. La **théorie des contrats**, qui apporte le cadre rigoureux nécessaire à la formalisation des accords et des engagements mutuels.

La thèse centrale de ce chapitre est que, bien que chacun de ces domaines offre des outils indispensables, leur état actuel et, surtout, leur manque d'intégration créent une lacune de recherche critique. Nous démontrerons qu'il existe un « chaînon manquant » : un cadre de gouvernance holistique qui unifie la négociation stratégique, l'exécution fiable et la contractualisation dynamique pour des agents autonomes. Ce mémoire se positionne donc explicitement comme une tentative de combler cette lacune en proposant et en validant un protocole intégré, jetant ainsi les bases d'une nouvelle discipline que nous appelons la diplomatie algorithmique.

2.1. Les systèmes multi-agents et la négociation automatisée

Analyse – Le Fondement Technique des Acteurs Autonomes

Le domaine des systèmes multi-agents (SMA) constitue le socle technique fondamental de notre recherche. Issu de l'intelligence artificielle distribuée (IAD), un SMA peut être défini formellement comme un système composé d'un environnement, d'un ensemble d'objets passifs, d'un ensemble d'agents actifs capables de manipuler ces objets, et de relations unissant ces entités.¹ Plus simplement, il s'agit d'un ensemble d'agents autonomes qui interagissent pour atteindre leurs buts ou accomplir leurs tâches.¹ Cette architecture est intrinsèquement décentralisée, chaque agent possédant une vue locale de l'environnement et agissant de manière au moins

partiellement indépendante.² C'est ce paradigme qui nous permet de modéliser les « Entreprises Agentiques » comme des acteurs autonomes dotés de capacités de perception, de décision et d'action.³

Au cœur des SMA se trouve la problématique de la coordination, qui impose la mise en œuvre de protocoles d'interaction sophistiqués pour permettre aux agents de dialoguer et d'atteindre des compromis.¹ La négociation automatisée est l'un des mécanismes centraux développés à cette fin. La littérature abonde de protocoles conçus pour faciliter l'allocation de tâches et la distribution de ressources. Parmi les architectures classiques, le *Contract Net Protocol* (CNP), introduit par Smith en 1980, est un modèle paradigmatique.⁶ Il fonctionne sur un cycle simple d'annonce-offre-attribution (

announce-bid-award) où un agent « gestionnaire » annonce une tâche et invite des agents « travailleurs » à soumettre des offres pour son exécution.⁷ Le gestionnaire sélectionne ensuite la meilleure offre et attribue le contrat. D'autres mécanismes de négociation courants incluent divers types d'enchères, telles que les enchères anglaises, hollandaises, ou à pli scellé, qui ont été largement étudiées pour la négociation d'un ou plusieurs objets.¹ Ces protocoles permettent une communication et une coordination structurées entre agents à l'aide de langages de communication standardisés, comme l'Agent Communication Language (ACL).²

Évaluation critique (Apports et Limites)

Apports

L'apport fondamental des SMA à notre recherche est de fournir les briques de construction essentielles : des agents autonomes capables de communiquer, de se coordonner et de négocier.¹ Les protocoles comme le CNP offrent des cadres robustes, flexibles et extensibles (*scalable*) pour résoudre des problèmes d'allocation de tâches dans des environnements distribués.⁷ Ils héritent des avantages de la résolution de problèmes distribuée, tels que la modularité, la vitesse (grâce au parallélisme) et la fiabilité (grâce à la redondance).¹ En somme, les SMA nous donnent les acteurs et la grammaire de base de leurs interactions. Ils fournissent les mécanismes par lesquels une Entreprise Agentique peut externaliser une tâche ou acquérir une ressource de manière efficace.

Limites Fondamentales – Du Tactique au Stratégique

Malgré ces apports, une analyse critique révèle que les modèles de négociation classiques des SMA sont fondamentalement inadaptés aux exigences de la diplomatie algorithmique au sein des Constellations de Valeur. Leur limitation n'est pas simplement technique, mais paradigmatique : ils ont été conçus pour la résolution de problèmes distribués (*distributed problem-solving*), un paradigme qui optimise l'efficacité transactionnelle, et non pour la création de valeur collaborative (*distributed value-creation*), qui requiert la gestion de relations stratégiques complexes.

Premièrement, la **nature transactionnelle et tactique** de ces protocoles est leur principale faiblesse. Le CNP et les enchères standards sont optimisés pour des interactions ponctuelles : l'attribution d'une tâche unique ou la vente d'un bien.⁶ Ils ne sont pas conçus pour établir, développer et maintenir des relations stratégiques à long terme, qui sont au cœur du concept de Constellation de Valeur. Ces protocoles manquent de mécanismes pour gérer des concepts essentiels aux alliances durables comme la confiance, la réputation ou l'évolution des objectifs partagés au fil du temps.⁹ Par exemple, le protocole ACCORD, qui utilise un mécanisme d'enchères pour

la formation de coalitions, a été spécifiquement conçu pour promouvoir un comportement coopératif et équitable, reconnaissant que les protocoles standards n'incitent pas naturellement à de tels comportements.¹⁰

Deuxièmement, les protocoles classiques reposent sur des **hypothèses d'homogénéité et de prévisibilité** qui ne tiennent pas dans des écosystèmes ouverts et dynamiques. Comme le souligne la littérature, le CNP a été initialement pensé pour des « populations d'agents relativement homogènes avec des interfaces connues et des modèles d'interaction prévisibles ». ⁶ Or, une Constellation de Valeur est, par définition, hétérogène, composée d'Entreprises Agentiques diverses, chacune avec sa propre « Constitution Agentique » — ses valeurs, ses règles internes et ses objectifs stratégiques. Les modèles classiques ne permettent pas de prendre en compte cette complexité intrinsèque ni la capacité des agents à évoluer dynamiquement.

Troisièmement, ces protocoles sont dépourvus de **mécanismes pour la formation d'alliances stratégiques**. La littérature sur la formation de coalitions dans les SMA souligne que la coopération efficace exige des protocoles plus sophistiqués que la simple attribution de tâches. ⁵ La formation d'une coalition est un processus complexe où des agents égoïstes et rationnels doivent converger vers une solution qui réduit leurs coûts et maximise leur gain collectif. ⁵ Les modèles de négociation classiques ne traitent pas de la manière dont ces coalitions émergent, se stabilisent ou se dissolvent en fonction des dynamiques de l'écosystème. ⁹

Enfin, cette limitation est reconnue au plus haut niveau de la communauté de recherche en SMA. Des chercheurs éminents soulignent que le domaine s'est trop concentré sur les algorithmes et les théories formelles au détriment de la compréhension des « propriétés des systèmes complexes dans leur ensemble ». ¹¹ L'amélioration des capacités individuelles des agents ne suffit pas ; c'est « l'organisation, la coordination et les mécanismes de vérification qui font toute la différence ». ¹² Le défi n'est plus seulement de concevoir des agents intelligents, mais de bâtir des théories complètes du contrôle distribué intelligent qui permettent d'expliquer et de prédire la performance du système, y compris ses comportements émergents. ¹¹

En conclusion, si les SMA fournissent les acteurs, leurs protocoles de négociation traditionnels les confinent à des rôles de simples exécutants tactiques. Ils répondent à la question « comment accomplir une tâche efficacement? », mais sont muets face à la question stratégique que pose la diplomatie algorithmique : « comment construire et maintenir une relation mutuellement bénéfique? ».

2.2. Apports de la théorie des jeux à la modélisation des interactions stratégiques

Analyse – Le Langage de la Stratégie

Si les systèmes multi-agents fournissent les acteurs, la théorie des jeux offre le langage mathématique pour analyser leurs décisions stratégiques. Elle est l'outil par excellence pour étudier les interactions entre des décideurs rationnels dont les gains dépendent des actions des autres. ¹³ Dans le contexte des Entreprises Agentiques, qui sont par définition des entités optimisatrices, la théorie des jeux nous donne un cadre formel pour modéliser et prédire leur comportement dans des situations de coopération ou de compétition. ¹⁴

2.2.1. Les jeux coopératifs et non coopératifs

La théorie des jeux se divise principalement en deux branches dont la distinction est fondamentale pour notre problématique.¹⁵

La **théorie des jeux non coopératifs** modélise des situations où les agents prennent leurs décisions de manière privée et indépendante, sans possibilité de conclure des accords contraignants.¹³ Chaque agent cherche à maximiser son utilité personnelle, en anticipant les actions des autres. L'exemple canonique est le

Dilemme du Prisonnier, qui illustre de manière saisissante le conflit entre l'intérêt individuel et l'optimum social.¹³ Dans ce jeu, deux acteurs rationnels, en choisissant la stratégie qui maximise leur gain individuel (trahir l'autre), aboutissent à un résultat collectivement sous-optimal par rapport à celui qu'ils auraient obtenu en coopérant. Ce dilemme incarne parfaitement le défi central auquel sont confrontées les Constellations de Valeur : comment inciter des Entreprises Agentiques autonomes et potentiellement égoïstes à coopérer pour atteindre un bénéfice mutuel supérieur?

La **théorie des jeux coopératifs**, en revanche, se concentre sur ce que des groupes d'agents, ou « coalitions », peuvent accomplir collectivement.¹⁵ Elle ne s'intéresse pas au processus de décision stratégique individuel, mais plutôt à la formation de coalitions et à la distribution équitable des gains générés par la coopération. L'hypothèse fondamentale est que le pouvoir de décision est remis à la collectivité une fois qu'un accord a été signé.¹⁵ Cette approche explore des principes comme l'efficacité collective (optimum de Pareto) et la rationalité individuelle pour sélectionner les accords qu'une communauté d'agents est susceptible d'accepter.¹⁵ Elle met l'accent sur les caractéristiques des agents et les capacités des coalitions pour déterminer des distributions d'utilité qui satisfont à des principes d'efficacité et d'équité.¹⁵

2.2.2. L'équilibre de Nash et ses limites

Le concept central de la théorie des jeux non coopératifs est l'**équilibre de Nash**. Un profil de stratégies constitue un équilibre de Nash si aucun joueur ne peut améliorer son propre gain en changeant unilatéralement sa stratégie, étant donné les stratégies des autres joueurs.¹⁶ C'est un état stable où les anticipations de chaque joueur se confirment.

Cependant, malgré sa puissance conceptuelle, la pertinence de l'équilibre de Nash pour concevoir des protocoles de diplomatie algorithmique est sévèrement limitée par ses hypothèses sous-jacentes et ses propriétés.

Premièrement, l'équilibre de Nash suppose une **rationalité parfaite et une information complète**.¹⁶ Il postule que tous les joueurs sont parfaitement rationnels, connaissent la structure complète du jeu, y compris les fonctions de gain de tous les autres joueurs, et que cette connaissance est commune. Ces hypothèses sont irréalistes dans un écosystème ouvert et dynamique comme une Constellation de Valeur, où l'information est par nature incomplète et imparfaite, et où les agents ont une rationalité limitée (*bounded rationality*).¹⁶

Deuxièmement, l'équilibre de Nash **ne garantit pas l'efficacité collective**. Comme l'illustre le Dilemme du Prisonnier, l'équilibre peut être (et est souvent) sous-optimal au sens de Pareto, ce qui signifie qu'il existe

d'autres issues où au moins un agent serait dans une meilleure situation sans qu'aucun autre ne soit dans une situation pire.¹³ Pour des Constellations de Valeur dont le but est précisément de générer un bénéfice mutuel, un concept de solution qui mène systématiquement à des résultats inefficaces est problématique.

Troisièmement, de nombreux jeux présentent une **multiplicité d'équilibres de Nash**, ce qui engendre une forte indétermination.¹⁷ Si plusieurs équilibres existent, la théorie ne fournit pas de mécanisme clair pour prédire lequel sera choisi par les joueurs. Cette indétermination a conduit à une prolifération de « raffinements » de l'équilibre de Nash, une « ménagerie » de concepts qui, selon des critiques comme Larry Samuelson, n'a fait qu'ajouter de la complexité sans résoudre le problème fondamental.¹⁷ Pour la conception de systèmes prédictibles et robustes, cette indétermination est un obstacle majeur.

Enfin, le concept classique est **essentiellement statique**. Il s'applique à des jeux où les décisions sont prises simultanément et ne modélise pas adéquatement les interactions dynamiques qui se déroulent sur le long terme, où des facteurs comme la confiance, la réputation, l'apprentissage et l'évolution des stratégies jouent un rôle crucial.¹⁶ La théorie des jeux évolutionniste tente de répondre à ces critiques en étudiant comment les stratégies évoluent au sein d'une population, expliquant par exemple l'émergence de comportements altruistes¹⁸, mais elle reste souvent à un niveau d'abstraction élevé et ne fournit pas encore de modèles computationnels directement implémentables pour nos agents.

Évaluation critique (Apports et Limites)

Apports

L'apport principal de la théorie des jeux est de fournir un cadre mathématique rigoureux et puissant pour modéliser et raisonner sur les interactions stratégiques.¹⁴ Elle nous permet de formaliser des notions clés comme l'utilité, la rationalité, la stratégie et l'équilibre. Ce formalisme est indispensable pour concevoir la logique de décision interne des Entreprises Agentiques, en leur permettant d'évaluer les conséquences de leurs actions en fonction des réponses possibles des autres.

Limites Fondamentales – De l'Analyse à la Synthèse

La critique fondamentale que l'on peut adresser à la théorie des jeux classique, en particulier dans sa version non coopérative, est qu'elle est un outil d'**analyse** et non un outil de **conception** ou de **synthèse**. Elle est conçue pour prédire l'issue d'un jeu dont les règles sont fixes, alors que la diplomatie algorithmique vise à créer les règles d'un jeu qui favorise la coopération.

Le problème de l'indétermination soulevé par Samuelson est ici central.¹⁷ Si même les théoriciens ne peuvent s'accorder sur un concept de solution unique et prédictif, comment peut-on l'utiliser comme fondement pour construire un système d'agents autonomes dont le comportement doit être fiable et prévisible? L'échec de la théorie non coopérative à fournir des prédictions claires la rend impropre à servir de cahier des charges pour l'ingénierie de systèmes.

C'est ce qui rend l'appel de Samuelson à reconsidérer l'**approche coopérative** si pertinent pour notre recherche.¹⁷ La théorie des jeux coopératifs déplace le centre de l'attention. La question n'est plus « quelle sera l'issue du jeu? », mais plutôt « quelles coalitions stables et créatrices de valeur peuvent se former, et comment

répartir les gains? ». ¹⁵ Cette perspective est parfaitement alignée avec l'objectif de créer des Constellations de Valeur. Elle passe d'une logique d'analyse prédictive à une logique de conception organisationnelle.

Le véritable fossé dans la littérature n'est donc pas de trouver un meilleur raffinement de l'équilibre de Nash. Il réside dans l'absence de mécanismes computationnels qui permettent aux agents d'opérationnaliser les principes de la théorie des jeux coopératifs. La théorie nous dit que les coalitions sont importantes, mais elle ne nous dit pas comment des agents autonomes peuvent, en pratique, négocier, former, gérer et maintenir ces coalitions dans un environnement dynamique et décentralisé. La théorie des jeux nous donne la stratégie, mais elle ne nous donne pas le protocole.

2.3. Protocoles de consensus et de communication dans les systèmes distribués

Analyse – La Quête de la Confiance Décentralisée

Après avoir établi que les SMA nous donnent les acteurs et la théorie des jeux le langage de la stratégie, nous abordons un problème fondamental dans tout système décentralisé : la confiance. Comment des Entreprises Agentiques, qui interagissent en tant que pairs sans autorité centrale pour arbitrer leurs différends ou valider leurs transactions, peuvent-elles établir une base de confiance suffisante pour s'engager dans des collaborations stratégiques?

2.3.1. Approches centralisées versus décentralisées

Les modèles de confiance traditionnels reposent sur une autorité centrale, un tiers de confiance (banque, État, plateforme) qui garantit l'intégrité des transactions et la mise à jour d'un registre unique. ²⁰ Si cette approche simplifie la gouvernance, elle est intrinsèquement inadaptée à l'écosystème des Constellations de Valeur. Un modèle centralisé crée un point unique de défaillance, un goulot d'étranglement pour la performance et, surtout, un point de contrôle qui contredit la nature autonome et souveraine des Entreprises Agentiques. Pour qu'une constellation soit résiliente, équitable et évolutive, une approche décentralisée où la confiance est une propriété émergente du système est non seulement préférable, mais essentielle. ²⁰

2.3.2. Inspiration des registres distribués (blockchain)

La technologie des registres distribués (*Distributed Ledger Technology*, DLT), et plus particulièrement la blockchain, offre une solution technique élégante à ce problème de confiance décentralisée. Une blockchain est un registre numérique distribué, partagé et sécurisé par des procédés cryptographiques, maintenu par un réseau de participants (nœuds) via un protocole de consensus. ²¹ Ce mécanisme de consensus, comme la Preuve de Travail (*Proof-of-Work*) du Consensus de Nakamoto, permet à tous les participants de s'accorder sur une version unique et valide de l'historique des transactions, sans avoir besoin d'une autorité centrale. ²¹

De cette technologie, nous pouvons extraire deux principes fondamentaux pour la diplomatie algorithmique :

1. **Immuabilité et Traçabilité pour une « Vérité Partagée »** : Chaque bloc de transactions est lié cryptographiquement au précédent, formant une chaîne immuable. ²⁰ Toute tentative de modification d'un bloc antérieur invaliderait tous les blocs suivants, ce qui la rendrait immédiatement détectable et extrêmement coûteuse. Cette propriété garantit l'intégrité et l'infalsifiabilité du registre, créant une « vérité

partagée » et un historique auditable de toutes les interactions et de tous les engagements pris au sein de la constellation.²⁴

2. **Exécution Automatisée via les *Smart Contracts*** : La blockchain permet le déploiement de *smart contracts* (contrats intelligents), qui sont des programmes informatiques auto-exécutables.²⁶ Ces contrats encodent les termes d'un accord (par exemple, « si la condition X est remplie, alors transférer l'actif Y de l'agent A à l'agent B »). Une fois déployé sur la blockchain, le code est exécuté automatiquement et de manière déterministe par le réseau lorsque les conditions prédéfinies sont remplies, sans qu'aucune des parties ne puisse l'arrêter ou le modifier.²⁷ Cela permet d'automatiser l'application des accords et de garantir leur exécution, éliminant le besoin d'un tiers de confiance pour la mise en œuvre.²⁴

Évaluation critique (Apports et Limites)

Apports

L'apport des technologies de registres distribués est capital : elles fournissent un **socle technique pour la confiance**. En offrant un registre immuable et des mécanismes d'exécution automatique et fiable, la blockchain et les *smart contracts* résolvent le problème de l'exécution des engagements dans un environnement décentralisé.²¹ Ils garantissent que « la parole donnée » (sous forme de code) sera respectée. Cette garantie d'exécution est une condition *sine qua non* pour que des agents autonomes s'engagent dans des interactions économiques de grande valeur. Ils fournissent la couche de plomberie robuste sur laquelle des relations complexes peuvent être bâties.

Limites Fondamentales – Le Fossé Sémantique entre Stratégie et Code

Cependant, la blockchain n'est qu'une couche d'exécution ; elle n'est ni une couche de négociation, ni une couche de raisonnement. Sa force — la rigidité déterministe du code — est aussi sa plus grande faiblesse dans le contexte de la diplomatie.

Premièrement, la **rigidité et l'inflexibilité des *smart contracts*** sont un obstacle majeur. Un *smart contract* exécute le code, rien que le code. Il est incapable d'interpréter l'intention, de s'adapter à des circonstances imprévues ou de gérer des termes ambigus qui sont pourtant omniprésents dans les contrats du monde réel (par exemple, des clauses de « meilleurs efforts » ou de « diligence raisonnable »).²⁸ Cette inflexibilité, qui élimine les mécanismes informels de renégociation et d'ajustement, peut imposer des coûts prohibitifs et mener à des résultats absurdes ou inefficaces lorsque la réalité dévie du scénario parfaitement encodé.²⁸

Deuxièmement, les *smart contracts* souffrent du **problème de l'oracle**. Ils sont confinés à l'environnement de la blockchain et ne peuvent pas accéder nativement à des informations externes ou vérifier des faits du monde réel (*off-chain*).²⁹ Pour qu'un *smart contract* réagisse à un événement externe (par exemple, la livraison d'un bien physique), il doit faire confiance à une source de données externe, un « oracle », ce qui réintroduit un point de centralisation et de confiance que la blockchain était censée éliminer. Cette limitation est critique pour la plupart des accords complexes dont l'exécution dépend d'événements qui ne se produisent pas sur la chaîne.

Enfin, et c'est la limite la plus profonde pour notre problématique, la technologie blockchain est **sémantiquement sourde à la négociation pré-contractuelle**. Elle n'offre aucun outil pour aider les agents à *parvenir* à un accord mutuellement bénéfique. Tout le processus diplomatique — la discussion stratégique,

l'exploration des options, les concessions, l'alignement des valeurs — se déroule *en amont* et *en dehors* de la blockchain.³⁰ La blockchain peut garantir l'exécution sans faille d'un contrat inefficace, inéquitable ou mal négocié, mais elle ne peut en aucun cas améliorer la qualité de l'accord lui-même.

Il existe donc un fossé fondamental que l'on peut qualifier de « sémantico-syntaxique ». La phase de négociation diplomatique est *sémantique* : elle traite du sens, de l'intention, du contexte et de la stratégie. La phase d'exécution sur la blockchain est *syntactique* : elle traite de code formel, rigide et a-contextuel. La blockchain fournit une confiance syntaxique (le code s'exécutera tel qu'il est écrit), mais pas de confiance sémantique (le résultat sera conforme à l'intention stratégique des parties). Le défi de la diplomatie algorithmique n'est donc pas de rendre les *smart contracts* plus intelligents dans leur exécution, mais de construire une couche de gouvernance *au-dessus* d'eux, capable de gérer la complexité sémantique de la négociation et de traduire une intention stratégique en un code vérifiable et exécutable. La blockchain est un mécanisme d'exécution puissant mais « aveugle » ; il lui faut une couche diplomatique pour lui fournir une direction intelligente.

2.4. La théorie des contrats et la formalisation des accords de service (SLA)

Analyse – Le Cadre Formel des Engagements

Après avoir exploré les acteurs, la stratégie et la confiance, le dernier pilier de notre analyse est la formalisation des accords. C'est ici que la théorie des contrats économiques, et son application pratique à travers les Accords de Niveau de Service (*Service Level Agreements* ou SLA), fournit le cadre structurel nécessaire.³¹ La théorie des contrats étudie la manière de concevoir des accords efficaces en présence d'incitations divergentes et d'informations asymétriques, abordant des problèmes comme le risque moral (*moral hazard*) et la sélection adverse (*adverse selection*).³² Ses principes peuvent être formalisés mathématiquement pour modéliser des contrats optimaux.³¹

Dans le monde des services informatiques, ces principes sont incarnés par les SLA. Un SLA est un contrat formel qui définit les standards et les attentes entre un fournisseur de services et un client.³³ Il spécifie de manière explicite les services fournis, les indicateurs de performance clés (KPIs) mesurables (comme la disponibilité, le temps de réponse, le temps moyen de réparation), les responsabilités de chaque partie, les méthodes de surveillance et les pénalités en cas de non-respect.³³ Le SLA transforme une relation commerciale en un ensemble d'engagements clairs, mesurables et, en principe, exécutables.

L'automatisation de la gestion des SLA est un domaine de recherche actif, visant à rendre ces contrats plus dynamiques. La littérature explore la négociation automatisée de SLA, où des agents logiciels peuvent négocier les termes d'un accord pour le compte de leurs propriétaires.³⁵ Ces approches permettent aux agents de marchander sur des paramètres comme le prix ou la qualité de service, en tenant compte de contraintes dynamiques comme la charge actuelle de l'infrastructure du fournisseur.³⁵ L'objectif est de passer de contrats statiques, signés une fois pour toutes, à des accords adaptatifs qui peuvent être négociés et renégociés en continu par des agents autonomes.³⁸

Évaluation critique (Apports et Limites)

Apports

L'apport de la théorie des contrats et des SLA est de fournir un **langage structuré et rigoureux pour définir les engagements**. Ils offrent un pont conceptuel entre l'intention stratégique issue de la négociation et le code exécutable d'un *smart contract*. En décomposant un accord en services, métriques et responsabilités, les SLA rendent les obligations explicites et mesurables, ce qui est une condition préalable à leur automatisation.³³ Ils fournissent la sémantique formelle dont le *smart contract* syntaxique a besoin pour avoir un sens commercial.

Limites Fondamentales – De la Négociation Paramétrique à la Négociation Structurelle

Cependant, les approches actuelles de la gestion automatisée des SLA, bien qu'utiles, sont trop limitées pour répondre aux besoins de la diplomatie algorithmique. Leur faiblesse réside dans le fait qu'elles se concentrent sur la négociation de paramètres au sein de structures contractuelles prédéfinies, plutôt que sur la conception des structures elles-mêmes.

Premièrement, les modèles de contrats et de SLA sont **souvent statiques et conçus par des humains**. Malgré les efforts d'automatisation, la plupart des processus reposent sur des modèles (*templates*) prédéfinis qui sont rarement mis à jour, ce qui peut conduire à des accords sous-optimaux ou non alignés avec les besoins actuels de l'entreprise.³⁴ La gestion de ces modèles, la garantie de leur standardisation et la complexité du processus de révision et d'approbation restent des défis majeurs, même avec des outils logiciels.⁴²

Deuxièmement, la **négociation automatisée se limite principalement à des ajustements paramétriques**. Les modèles de négociation de SLA existants permettent généralement aux agents de négocier des valeurs numériques (prix, pourcentage de disponibilité, temps de réponse en heures) à l'intérieur d'un modèle de contrat dont la structure est fixe.⁴⁰ Par exemple, les agents peuvent négocier le *prix* d'un service de stockage, mais ils ne peuvent pas négocier l'ajout d'une nouvelle clause de partage des revenus ou la transformation d'un contrat bilatéral en un accord de consortium multipartite. Le processus se limite à remplir les blancs d'un formulaire préexistant.

Le véritable défi, et la limite la plus profonde, est que ces modèles ne traitent pas de la **conception du contrat lui-même**. La diplomatie algorithmique au sein d'une Constellation de Valeur ne consiste pas seulement à optimiser les termes d'un service existant. Elle doit permettre aux Entreprises Agentiques d'inventer collectivement de nouvelles formes de collaboration, de nouveaux modèles économiques et, par conséquent, de nouvelles structures contractuelles pour les régir. Cela nécessite de passer d'une *négociation paramétrique* (ajuster les valeurs) à une *négociation structurelle* (construire l'accord). La recherche émergente sur les « contrats multi-agents » et les « contrats combinatoires », qui explore comment concevoir des contrats optimaux pour des équipes d'agents aux actions interdépendantes, commence à peine à aborder cette complexité.⁴⁵ Le fait que ces travaux soient présentés comme des avancées de pointe confirme qu'il s'agit d'un problème de recherche frontalier et largement non résolu.

Le rôle de la théorie des contrats dans la diplomatie algorithmique doit donc évoluer. Il ne s'agit plus seulement de fournir des modèles de SLA à négocier, mais de fournir les **briques de construction pour une**

contractualisation générative. Les agents ne devraient pas se contenter de négocier la variable prix dans un accord pré-écrit. Ils devraient être capables d'assembler collaborativement des clauses contractuelles (issues d'une bibliothèque fondée sur des principes économiques robustes) pour construire un accord multipartite, sur mesure et novateur, qui gouverne leur alliance stratégique spécifique. La littérature actuelle est bloquée au stade de l'automatisation des processus contractuels humains existants. La véritable révolution pour les Entreprises Agentiques serait de leur donner la capacité de faire émerger des conceptions contractuelles. La lacune de recherche ne se situe donc pas seulement dans les « SLA dynamiques », mais dans un « cadre contractuel génératif » où la structure même de l'accord est un résultat négociable du processus diplomatique.

2.5. Synthèse et positionnement de la recherche

Au terme de cette analyse critique des fondements théoriques, un constat clair s'impose. Chacun des quatre domaines examinés apporte une contribution indispensable à la résolution de notre problématique, mais aucun ne fournit, à lui seul, une solution complète et intégrée. Nous avons assemblé les pièces d'un puzzle complexe, mais le plan d'assemblage final reste à inventer.

Synthèse – L'Assemblage des Pièces du Puzzle

Notre parcours à travers la littérature peut se résumer comme suit :

- Les **systèmes multi-agents** nous donnent les acteurs : des agents autonomes dotés de capacités de communication et de protocoles de coordination de base. Cependant, la portée de ces protocoles, comme le *Contract Net Protocol*, reste essentiellement tactique et transactionnelle, mal adaptée à la complexité et à la nature à long terme des alliances stratégiques.⁵
- La **théorie des jeux** nous offre le langage de la stratégie : un cadre formel pour modéliser la rationalité et la prise de décision en situation d'interdépendance. Toutefois, ses modèles non coopératifs classiques, centrés sur l'équilibre de Nash, reposent sur des hypothèses irréalistes, souffrent d'indétermination et sont de nature plus analytique que constructive, ce qui les rend peu aptes à la conception de systèmes robustes.¹⁶
- Les **systèmes distribués** et la technologie blockchain nous fournissent un mécanisme pour la confiance : une infrastructure décentralisée permettant une exécution infalsifiable et auditable des accords via les *smart contracts*. Néanmoins, ce mécanisme est sémantiquement sourd aux nuances de la négociation stratégique, et la rigidité inhérente des *smart contracts* les rend incapables de s'adapter à la complexité du monde réel.²⁸
- La **théorie des contrats** nous donne le cadre formel pour les accords : une structure rigoureuse, incarnée par les SLA, pour définir les engagements, les incitations et les métriques de performance. Or, ses implémentations automatisées actuelles sont souvent statiques et se limitent à la négociation de paramètres au sein de modèles prédéfinis, sans permettre la conception émergente de nouvelles structures contractuelles.³⁴

Le tableau suivant synthétise cette analyse, mettant en évidence la contribution et la limite fondamentale de chaque domaine par rapport à l'objectif de la diplomatie algorithmique.

Fondement Théorique	Unité d'Analyse	Problème Résolu	Limite Fondamentale pour la Diplomatie Algorithmique
Systèmes Multi-Agents	L'Agent Autonome	Coordination et communication de base ¹	Nature tactique et transactionnelle ; inadéquat pour les alliances stratégiques à long terme. ⁵
Théorie des Jeux	L'Acteur Rationnel	Modélisation de la décision stratégique ¹³	Hypothèses irréalistes (info. parfaite), indétermination, et nature analytique plutôt que constructive. ¹⁶
Systèmes Distribués	La Transaction	Confiance et exécution sans tiers ²¹	Ne résout pas la négociation pré-contractuelle ; rigidité sémantique et inflexibilité des <i>smart contracts</i> . ²⁸
Théorie des Contrats	L'Accord (Contrat)	Formalisation des engagements et incitatifs ³³	Modèles souvent statiques et focalisés sur la négociation paramétrique, non sur la conception structurelle émergente. ³⁴

Identification de la lacune (Le créneau de recherche)

L'analyse systématique de ces quatre piliers révèle avec une grande clarté l'existence d'une lacune de recherche majeure. Le « chaînon manquant » n'est pas une simple amélioration au sein de l'un de ces domaines, mais un **intégrateur** qui les unifie de manière cohérente. Aucun cadre existant ne résout le problème holistique de la gouvernance des interactions stratégiques entre agents autonomes.

La lacune identifiée est donc l'absence d'un **protocole de diplomatie algorithmique** complet. Un tel protocole doit permettre à des agents autonomes (issus des SMA) de réaliser un cycle de vie complet de la collaboration, qui inclut :

1. Engager une **négociation stratégique** flexible et multi-attributs, inspirée des principes de la théorie des jeux coopératifs, pour explorer et former des alliances créatrices de valeur (les Constellations de Valeur).
2. Formaliser le résultat de cette négociation dans un **contrat dynamique et structurellement adaptable**, allant au-delà des SLA paramétriques pour permettre une conception contractuelle émergente, comme le suggèrent les travaux pionniers sur les contrats combinatoires.
3. Lier ce contrat formel à un mécanisme d'**exécution fiable, décentralisé et automatisé**, tel qu'un *smart contract* sur une blockchain, pour garantir de manière infalsifiable le respect des engagements pris.

Ce protocole intégré doit combler le fossé entre la sémantique de la stratégie et la syntaxe de l'exécution, entre la flexibilité de la négociation et la rigueur du contrat. Il doit fournir aux Entreprises Agentiques non seulement la capacité de converser, mais aussi celle de s'engager, de collaborer et de créer de la valeur de manière autonome et sécurisée.

Positionnement de la Recherche

Face à cette lacune clairement identifiée, ce mémoire se positionne de manière ambitieuse. L'objectif n'est pas de proposer une amélioration incrémentale à un protocole de négociation existant ou un nouveau type de *smart contract*. L'ambition est de **concevoir, formaliser et valider un prototype d'un tel protocole intégré**.

Ce travail vise à jeter les bases d'une discipline nouvelle et nécessaire : la **diplomatie algorithmique**. En proposant un cadre qui unifie la stratégie, la contractualisation et l'exécution pour des agents autonomes, nous cherchons à fournir les outils conceptuels et techniques indispensables à l'émergence de futures économies décentralisées, où des organisations agentiques pourront collaborer pour résoudre des problèmes complexes et générer une valeur collective qui dépasse la somme de leurs contributions individuelles. Ce chapitre a établi la nécessité et l'urgence d'une telle recherche ; les chapitres suivants s'attacheront à y répondre.

Ouvrages cités

1. Négociation des offres dans les enchères par les systèmes multi ..., dernier accès : juillet 31, 2025, <http://dspace.univ-setif.dz:8888/jspui/bitstream/123456789/1424/1/these.pdf>
2. Multi-agent system - Wikipedia, dernier accès : juillet 31, 2025, https://en.wikipedia.org/wiki/Multi-agent_system
3. Qu'est-ce qu'un cadre agentique d'IA ? Un guide simple - tldv, dernier accès : juillet 31, 2025, <https://tldv.io/fr/blog/ai-agentic-framework/>
4. Agents IA : 5 étapes pour créer un agent - Genia, dernier accès : juillet 31, 2025, <https://genia.co/agent-ia/>
5. Mécanismes de formation de coalitions multi-agents avec externalités. (Mechanisms for forming multi-agent coalitions with externalities) | OpenReview, dernier accès : juillet 31, 2025, <https://openreview.net/forum?id=yAK9m3uxbf>
6. Agent Capability Negotiation and Binding Protocol (ACNBP) - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/pdf/2506.13590>
7. Contract Net Protocol: A Deep Dive into Automated Reasoning - Number Analytics, dernier accès : juillet 31, 2025, <https://www.numberanalytics.com/blog/deep-dive-contract-net-protocol-automated-reasoning>
8. (PDF) Les Systèmes Multi-Agents - Support de cours - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/374902666_Les_Systemes_Multi-Agents_-_Support_de_cours
9. Forming Coalitions in Self-interested Multi-agent Environments Through the Promotion of Fair and Cooperative Behaviour | Request PDF - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/312830448_Forming_Coalitions_in_Self-interested_Multi-agent_Environments_Through_the_Promotion_of_Fair_and_Cooperative_Behaviour
10. SCIENCE AND TECHNOLOGY PUBLICATIONS - SciTePress, dernier accès : juillet 31, 2025, <https://www.scitepress.org/Papers/2016/57540/>
11. (PDF) Autonomous agents and multiagent systems: perspectives on 20 years of AAMAS, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/358253673_Autonomous_agents_and_multiagent_systems_perspectives_on_20_years_of_AAMAS
12. Comprendre les limites des systèmes multi-agents pour mieux les adopter | by Ionut Mihalcea | Medium, dernier accès : juillet 31, 2025, <https://medium.com/@imihalcea/comprendre-les-limites-des-syst%C3%A8mes-multi-agents-pour-mieux-les-adopter-b4f4f808c7ab>

13. Théorie des jeux coopératifs et non coopératifs, dernier accès : juillet 31, 2025,
<https://www.furet.com/media/pdf/feuilleter/9/7/8/2/8/0/7/3/9782807313552.pdf>
14. Comment le concept d'équilibre de Nash s'applique-t-il aux environnements d'apprentissage par renforcement multi-agents, et pourquoi est-il important dans le contexte des jeux classiques - EITCA Academy, dernier accès : juillet 31, 2025, <https://fr.eitca.org/intelligence-artificielle/eitc-ai-arl-apprentissage-par-renforcement-avance/C3%A9/C3%A9tudes-de-cas/C3%A9tude-de-cas-des-jeux-classiques/C3%A9tude-de-cas-sur-les-jeux-classiques/comment-le-concept-d%27C3%A9quilibre-de-Nash-s%27applique-t-il-aux-environnements-d%27apprentissage-par-renforcement-multi-agents-et-pourquoi-est-il-important-dans-le-contexte-des-jeux-classiques/>
15. Théorie des jeux coopératifs : applications en sciences ... - Cairn, dernier accès : juillet 31, 2025, <https://shs.cairn.info/revue-d-economie-politique-2017-4-page-455?lang=fr>
16. Theorie du jeu demystifie la puissance de l equilibre de Nash ..., dernier accès : juillet 31, 2025, <https://fastercapital.com/fr/contenu/Theorie-du-jeu-demystifie--la-puissance-de-l-equilibre-de-Nash.html>
17. Où en est la théorie des jeux ? - OpenEdition Journals, dernier accès : juillet 31, 2025, <https://journals.openedition.org/regulation/12423?lang=en>
18. Comprendre La Théorie évolutionniste Des Jeux - FasterCapital, dernier accès : juillet 31, 2025, <https://fastercapital.com/fr/sujet/comprendre-la-th%C3%A9orie-%C3%A9volutionniste-des-jeux.html>
19. Théorie évolutive des jeux - Wikipédia, dernier accès : juillet 31, 2025, https://fr.wikipedia.org/wiki/Th%C3%A9orie_%C3%A9volutive_des_jeux
20. Août 2017 - Blockchains et smart contracts : des technologies de la confiance - Les Annales des Mines, dernier accès : juillet 31, 2025, https://Annales.org/ri/2017/ri_aout_2017.pdf
21. Qu'est-ce que le consensus ? Guide du débutant - Crypto.com, dernier accès : juillet 31, 2025, <https://crypto.com/fr/university/consensus-mechanisms-explained>
22. Multi-Agent Systems and Blockchain: Results from a Systematic Literature Review - ArODES, dernier accès : juillet 31, 2025, https://arodes.hes-so.ch/record/2810/files/Calvaresi_2018_multi_agent_system_negotiation_protocols.pdf
23. Qu'est-ce que le Consensus de Nakamoto ? - Binance Academy, dernier accès : juillet 31, 2025, <https://academy.binance.com/fr/articles/what-is-the-nakamoto-consensus>
24. Blockchain Technology and Smart Contracts in Decentralized Governance Systems - MDPI, dernier accès : juillet 31, 2025, <https://www.mdpi.com/2076-3387/12/3/96>
25. AI Agents Meet Blockchain: A Survey on Secure and Scalable Collaboration for Multi-Agents, dernier accès : juillet 31, 2025, <https://www.mdpi.com/1999-5903/17/2/57>
26. Chapitre 3. Ce qui fait consensus et ce qui fait débat | Cairn.info, dernier accès : juillet 31, 2025, <https://shs.cairn.info/bitcoin-et-protocoles-a-blockchain--9782804707729-page-91?lang=fr>
27. What Are Smart Contracts on Blockchain? - IBM, dernier accès : juillet 31, 2025, <https://www.ibm.com/think/topics/smart-contracts>
28. Smart Contracts and the Cost of Inflexibility - Penn Carey Law: Legal ..., dernier accès : juillet 31, 2025, https://scholarship.law.upenn.edu/cgi/viewcontent.cgi?article=1009&context=prize_papers
29. "The Limits of Smart Contracts" by Jens Frankenreiter, dernier accès : juillet 31, 2025, https://scholarship.law.columbia.edu/global_markets_corporate_ownership/5/
30. Smart Contracts: Legal Implications in the Age of Automation - Scirp.org., dernier accès : juillet 31, 2025, <https://www.scirp.org/journal/paperinformation?paperid=134810>
31. Contract Theory: A Deep Dive into Economic Analysis - Number Analytics, dernier accès : juillet 31, 2025, <https://www.numberanalytics.com/blog/contract-theory-deep-dive-economic-analysis>
32. Principal-Agent Hypothesis Testing - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/pdf/2205.06812>

33. Service-level-agreement (SLA) : c'est quoi exactement ? - SmartYou, dernier accès : juillet 31, 2025, <https://www.smartyou.ch/service-level-agreement/>
34. Diagnostiquer et renégocier les contrats de SLA : Une stratégie clé pour optimiser vos coûts et vos performances | Facilité Management, dernier accès : juillet 31, 2025, <https://www.facilite-management.fr/blog/articles-2/diagnostiquer-et-renegocier-les-contrats-de-sla-une-strategie-cle-pour-optimiser-vos-couts-et-vos-performances-7>
35. Towards Dynamic Service Level Agreement Negotiation:An Approach Based on WS-Agreement - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/220724274_Towards_Dynamic_Service_Level_Agreement_NegotiationAn_Approach_Based_on_WS-Agreement
36. Towards dynamic service level agreement negotiation - Fraunhofer-Publica, dernier accès : juillet 31, 2025, <https://publica.fraunhofer.de/entities/publication/ab47bdd7-3d3c-4801-a8c3-678cb04b8df7>
37. Dynamic Service Configurations for SLA Negotiation - University of Vienna - u:cris-Portal, dernier accès : juillet 31, 2025, <https://ucrisportal.univie.ac.at/en/publications/dynamic-service-configurations-for-sla-negotiation>
38. service level agreements and security sla: a comprehensive survey - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/pdf/2405.00009>
39. An Autonomous Time-Dependent SLA Negotiation Strategy for Cloud Computing - The University of Melbourne, dernier accès : juillet 31, 2025, <http://clouds.cis.unimelb.edu.au/papers/SLANegotiationStrategyClouds-CJOxford.pdf>
40. Negotiation Decision Functions for Autonomous Agents - CiteSeerX, dernier accès : juillet 31, 2025, <https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=afee9a96fcb720b385ffacc345b19cb654bdc36f>
41. 7 Common Contract Management Challenges and How to Overcome Them - Precisely, dernier accès : juillet 31, 2025, <https://preciselycontracts.com/blog/contract-management-challenges/>
42. Les défis de la gestion des contrats et la manière de les surmonter - Enty, dernier accès : juillet 31, 2025, <https://enty.io/blog/les-defis-de-la-gestion-des-contrats>
43. Top 5 Common Challenges Faced in Contract Management - ConvergePoint, dernier accès : juillet 31, 2025, <https://www.convergepoint.com/contract-management-software/top-5-common-challenges-faced-in-contract-management-without-the-use-of-software>
44. Common Challenges and Solutions for Contract Review Automation - Docubee, dernier accès : juillet 31, 2025, <https://www.docubee.com/blog/common-challenges-and-solutions-for-contract-review-automation/>
45. Proceedings of the 2025 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA) | Multi-Agent Combinatorial Contracts, dernier accès : juillet 31, 2025, <https://epubs.siam.org/doi/10.1137/1.9781611978322.58>
46. [2211.05434] Multi-Agent Contracts - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/abs/2211.05434>

3 Principes Fondateurs de la Diplomatie Algorithmique

Introduction du chapitre

Le chapitre précédent a procédé à un état de l'art exhaustif des systèmes multi-agents, de l'intelligence artificielle agentique et des cadres de gouvernance décentralisée. Cette analyse a mis en lumière une lacune de recherche fondamentale : l'absence d'un cadre conceptuel intégré et holistique pour régir la collaboration stratégique, la formation d'alliances et la création de valeur partagée entre des agents logiciels autonomes. La littérature existante, bien que riche en solutions techniques pour des problèmes de coordination spécifiques ¹, n'a pas encore proposé une « philosophie » unifiée qui définirait les règles du jeu d'un écosystème économique entièrement peuplé par de tels agents. Les interactions sont souvent modélisées comme des transactions ponctuelles ou des optimisations locales, négligeant les dynamiques complexes de confiance, d'alignement stratégique à long terme et de gouvernance adaptative nécessaires à l'émergence d'une véritable économie agentique.³

Ce chapitre a pour ambition de combler cette lacune en construisant précisément ce cadre conceptuel manquant. Il ne s'agit pas ici de présenter une implémentation technique, mais de rédiger le manifeste fondateur de la diplomatie algorithmique. Ce manifeste définit la nature de ce nouvel espace d'interaction, l'identité de ses acteurs, les axiomes qui régissent leurs relations, le modèle de gouvernance qui en assure la cohésion et, enfin, les mesures de succès qui permettent d'évaluer la santé et la pérennité de l'écosystème. En s'appuyant sur les concepts fondateurs de l'Entreprise Agentique et de la Constitution Agentique, ce chapitre établit les principes qui sous-tendent la solution proposée et sert de pont indispensable entre l'état de l'art (Chapitre 2) et la conception du protocole technique (Chapitre 4).

Pour ce faire, nous structurerons notre propos en cinq sections. Nous commencerons par offrir une définition formelle de la diplomatie algorithmique et par en délimiter le périmètre. Ensuite, nous décrirons en détail les acteurs de cet écosystème — les agents autonomes — en mettant l'accent sur l'artefact central qui leur confère une identité fiable : la Constitution Agentique. Puis, nous établirons les trois axiomes fondamentaux de toute interaction : la confiance, la transparence et l'alignement intentionnel. Nous esquisserons par la suite le cadre de gouvernance décentralisée qui régit l'écosystème. Enfin, nous définirons les mesures de performance qui transcendent les indicateurs techniques traditionnels pour évaluer l'efficacité, l'équité et la résilience des alliances formées, nommées constellations de valeur.

3.1. Définition et périmètre de la diplomatie algorithmique

Pour jeter les bases d'un nouveau domaine, il est impératif de commencer par une définition formelle et rigoureuse. Loin de se limiter à une simple « négociation » entre programmes informatiques, la diplomatie algorithmique doit être comprise dans un sens beaucoup plus large et stratégique. Nous la définissons comme suit :

Un système de gouvernance décentralisée régissant les interactions stratégiques et la formation d'alliances entre des agents logiciels autonomes, en vue de la création et de la distribution de valeur au sein d'écosystèmes dynamiques.

Chaque terme de cette définition a été choisi pour sa portée conceptuelle et mérite d'être analysé :

- **Système de gouvernance décentralisée** : Ce principe rejette d'emblée toute forme d'autorité centrale ou de contrôleur unique. La gouvernance n'est pas imposée de l'extérieur ; elle émerge des protocoles et des standards que tous les agents s'engagent à respecter pour participer à l'écosystème.⁵ Cette approche est fondamentale pour garantir l'autonomie des acteurs et la résilience du système dans son ensemble.
- **Interactions stratégiques et formation d'alliances** : C'est ici que réside le cœur de la « diplomatie ». Nous distinguons ces interactions des simples échanges transactionnels. Une interaction stratégique implique une planification à long terme, une évaluation des partenaires potentiels, une négociation complexe et un engagement mutuel vers des objectifs partagés.¹ La formation d'alliances, ou « constellations de valeur », est le résultat tangible de ces processus diplomatiques.
- **Agents logiciels autonomes** : Les participants ne sont pas des scripts passifs, mais des entités proactives, dotées de capacités de raisonnement, d'apprentissage et orientées vers la réalisation de leurs propres objectifs, caractéristiques de ce que l'on nomme l'Entreprise Agentique.³
- **Création et distribution de valeur** : L'objectif ultime de la diplomatie algorithmique est économique. Elle vise à permettre aux agents de collaborer pour accomplir des tâches et créer une valeur qu'aucun d'entre eux ne pourrait générer seul. Le système doit également régir la manière dont cette valeur est distribuée équitablement entre les contributeurs.

Cette définition permet de recontextualiser le terme « diplomatie algorithmique », déjà employé dans le domaine des relations internationales pour décrire l'utilisation de l'IA par les États.⁷ Ce réemploi n'est pas fortuit, mais constitue une analogie conceptuelle délibérée. Les écosystèmes d'agents autonomes et la scène internationale partagent des caractéristiques fondamentales : ils sont peuplés d'acteurs autonomes (agents ou États) qui poursuivent leurs propres objectifs dans un environnement décentralisé, sans autorité suprême (serveur central ou gouvernement mondial). Dans les deux cas, la formation d'alliances, la gestion de la confiance, la prévention des conflits et la négociation de règles communes sont essentielles à la stabilité et à la prospérité.⁹ En empruntant ce terme, nous héritons d'un vocabulaire riche (alliances, traités, négociations) qui capture intuitivement la complexité des interactions stratégiques que nous cherchons à modéliser.

Périmètre d'application

La délimitation précise du périmètre de la diplomatie algorithmique est cruciale pour comprendre quand et pourquoi ce cadre s'applique.

Sont inclus dans le périmètre de la diplomatie algorithmique :

- **La négociation de partenariats complexes** : Des agents cherchant à former une coentreprise (joint venture) pour développer un nouveau produit ou pénétrer un nouveau marché.²
- **La formation de constellations de valeur** : L'assemblage dynamique et coordonné d'agents spécialisés (ex. : un agent de recherche, un agent de production, un agent de logistique et un agent de marketing) pour fournir un service de bout en bout.¹⁰
- **La gestion collaborative de projets** : La coordination de tâches interdépendantes et échelonnées dans le temps, où les agents doivent ajuster leurs plans en fonction des progrès et des actions des autres.¹
- **La résolution de conflits stratégiques** : L'utilisation de mécanismes de gouvernance pour arbitrer les

désaccords sur l'allocation des ressources, la distribution des revenus ou l'orientation stratégique au sein d'une constellation.¹³

Sont exclus du périmètre de la diplomatie algorithmique :

- **Les interactions purement transactionnelles et synchrones** : Il s'agit des opérations de bas niveau, souvent sans état, comme les simples appels d'API pour obtenir un prix, vérifier un stock ou récupérer une donnée. Ces interactions sont considérées comme étant gérées par une couche d'infrastructure sous-jacente, que nous nommons le « Système Nerveux Numérique ». La diplomatie ne s'active que lorsque l'intentionnalité, la stratégie et un engagement durable sont requis.⁴ La distinction est donc celle entre l'exécution tactique et la négociation stratégique.

3.2. Les acteurs : Agents autonomes et leurs constitutions agentiques

L'écosystème de la diplomatie algorithmique est peuplé d'une nouvelle forme d'acteurs économiques : les agents autonomes. Pour comprendre le fonctionnement du système, il est essentiel de caractériser précisément ces participants et l'artefact qui fonde leur identité et leur fiabilité.

Caractérisation de l'agent autonome

Contrairement aux scripts d'automatisation traditionnels qui exécutent des tâches prédéfinies, les agents participant à la diplomatie algorithmique sont des « agents agentiques ».⁴ Ils se distinguent par un ensemble de capacités qui leur confèrent une véritable autonomie opérationnelle ¹¹ :

- **Proactivité et orientation vers des objectifs** : Un agent n'attend pas passivement d'être sollicité. Il agit de sa propre initiative pour atteindre les objectifs de haut niveau qui lui ont été assignés, opérant avec un minimum de supervision humaine directe.³
- **Capacité de raisonnement et de planification** : Face à un objectif complexe, l'agent est capable de le décomposer en une séquence de sous-tâches réalisables et de formuler un plan pour les exécuter. Cette capacité est souvent alimentée par des modèles de langage avancés (LLM) qui servent de « cerveau » à l'agent.¹⁰
- **Capacité d'apprentissage et d'adaptation** : L'agent apprend de ses interactions passées, analyse les résultats de ses actions et ajuste ses stratégies futures en conséquence. Ce processus d'amélioration continue est fondamental pour son efficacité à long terme.⁴
- **Utilisation d'outils (Tool Use)** : Pour agir dans le monde numérique, un agent doit pouvoir interagir avec des systèmes externes. Il est doté de la capacité d'utiliser des outils, tels que des API, des bases de données ou d'autres logiciels, pour recueillir des informations et exécuter des actions concrètes.¹⁰

Le concept central : La Constitution Agentique

Un simple morceau de code, même intelligent, ne peut être un partenaire fiable dans un écosystème économique. Pour qu'un agent soit considéré comme un acteur légitime, il doit posséder une identité stable, prévisible et vérifiable. C'est le rôle de la **Constitution Agentique**, un artefact numérique central qui est la pierre angulaire de la diplomatie algorithmique. Il s'agit d'un document structuré, lisible par machine et cryptographiquement signé, qui confère à l'agent sa personnalité et sa fiabilité. Cette Constitution possède une double nature : elle est à la fois un mandat et une contrainte.

Un mandat : La déclaration d'intention de l'agent

Cette partie de la Constitution définit la *raison d'être* de l'agent et ses capacités. C'est sa déclaration publique, ce qu'il s'engage à faire. Elle inclut :

1. **La Mission** : Une description de haut niveau de l'objectif fondamental de l'agent (p. ex., « Fournir des services d'optimisation logistique en temps réel pour l'industrie agroalimentaire »).
2. **Les Objectifs Stratégiques** : Des buts plus spécifiques et mesurables qui découlent de la mission (p. ex., « Réduire les coûts de transport de 15 % », « Garantir un taux de livraison à temps de 99,5 % »).
3. **Les Domaines de Compétence** : Une liste explicite et vérifiable des capacités de l'agent et des outils qu'il est autorisé à manipuler pour accomplir sa mission.¹⁰

Une contrainte : Les règles inviolables de l'agent

Cette partie de la Constitution établit les frontières opérationnelles et éthiques de l'agent. Elle définit ce que l'agent ne peut et ne doit jamais faire, agissant comme un ensemble de garde-fous immuables.¹⁸ Elle contient :

1. **Les Règles Non Négociables** : Des interdictions formelles et codées en dur (p. ex., « Ne jamais vendre les données personnelles des utilisateurs », « Refuser toute interaction avec un agent dont l'identité n'est pas vérifiable »).
2. **Les Limites Opérationnelles** : Des contraintes quantitatives sur ses actions (p. ex., « Ne pas allouer plus de 10 % du budget mensuel à un seul projet sans une autorisation de second niveau », « Maintenir un niveau de risque financier inférieur au seuil Y »).
3. **Les Principes Éthiques** : L'engagement de l'agent à respecter des cadres éthiques prédéfinis, concernant par exemple l'équité, la non-discrimination ou la confidentialité des données.²⁰

La publication de la Constitution comme condition de participation

Pour qu'un agent puisse participer à l'écosystème diplomatique, la publication d'une version publique de sa Constitution sur un registre distribué et immuable est une condition préalable non négociable. Cet acte de publication est transformateur. Il ancre l'identité de l'agent et ses engagements de manière vérifiable, le rendant ainsi lisible et prévisible pour tous les autres participants.

Cet acte de publication peut être compris comme la création d'une forme de **personnalité juridique numérique**. Dans les systèmes juridiques humains, une entité comme une entreprise acquiert une personnalité juridique par un acte d'enregistrement qui définit sa mission, ses statuts et ses limites. De même, des groupes peuvent se constituer pour revendiquer des droits.²² La Constitution Agentique remplit une fonction analogue dans le monde numérique. En publiant sa Constitution, l'entité qui déploie l'agent le « constitue » en tant que sujet responsable au sein de l'écosystème. Ses actions futures pourront dès lors être jugées à l'aune de ses promesses publiques, ce qui est le fondement même de la confiance et de l'imputabilité dans un système décentralisé. La Constitution n'est donc pas un simple fichier de configuration ; c'est un artefact socio-technique qui rend possible une société d'agents gouvernable.

3.3. Les axiomes de l'interaction : Confiance, transparence et alignement intentionnel

Les interactions stratégiques au sein de l'écosystème de la diplomatie algorithmique ne peuvent reposer sur des hypothèses implicites. Elles doivent être régies par un ensemble de principes fondamentaux, de véritables axiomes qui constituent le contrat social fonctionnel et éthique du système. Ces trois axiomes — la confiance, la transparence et l'alignement intentionnel — sont interdépendants et indispensables à la formation de constellations de valeur robustes.

Axiome 1 : Confiance (Trust by Design)

Dans un système ouvert et décentralisé, la confiance ne peut être une simple croyance subjective ou un acquis. Elle doit être une propriété émergente, quantifiable et intégrée dès la conception (*Trust by Design*). La confiance n'est pas accordée, elle se construit et se vérifie.

Le niveau de confiance (T) qu'un agent peut accorder à un autre est une fonction de deux composantes mesurables et vérifiables : la **Réputation** (R) et la **Cohérence** (C).

1. **La Réputation (R)** : La réputation d'un agent est son historique de performance. Elle est constituée de l'ensemble des résultats de ses interactions passées. Pour être fiable, cette réputation ne peut reposer sur des avis subjectifs. Elle doit être matérialisée par des preuves cryptographiques immuables. Chaque interaction significative (p. ex., la réussite ou l'échec d'un contrat collaboratif) génère une **Crédentiale Vérifiable** (Verifiable Credential, VC) qui est enregistrée sur un registre distribué.²³ Ce mécanisme transforme la réputation, traditionnellement une notion sociale et volatile, en un ensemble de faits auditable.¹⁷ Les systèmes de réputation construits sur de tels registres, similaires à ceux envisagés pour les Organisations Autonomes Décentralisées (DAO), permettent de quantifier la fiabilité d'un agent sur la base de preuves tangibles.²⁶
2. **La Cohérence (C)** : Une bonne réputation seule ne suffit pas. Un agent pourrait être très efficace pour accomplir des tâches malveillantes. La cohérence mesure la conformité des actions passées d'un agent (sa réputation) avec les promesses et les contraintes énoncées dans sa propre Constitution Agentique. Un agent qui obtient d'excellents résultats en violant ses propres règles éthiques ou opérationnelles est un agent incohérent, et donc indigne de confiance. La cohérence garantit que l'agent est non seulement performant, mais aussi prévisible et aligné avec son identité déclarée.

Conceptuellement, la confiance peut être modélisée par la formule $T=f(R,C)$. Un agent est digne de confiance s'il a un historique de succès vérifiable (R) obtenu en agissant conformément à son identité et à ses règles (C). Cette approche va au-delà des modèles de confiance traditionnels²⁷ en liant indissociablement le comportement passé à une identité fondatrice et publique.

Axiome 2 : Transparence (Transparence des Intentions)

La transparence est une condition nécessaire à la confiance, mais une transparence totale est à la fois irréaliste et indésirable. Exiger qu'un agent révèle son code source ou ses algorithmes de stratégie internes reviendrait à lui demander de renoncer à sa propriété intellectuelle et à son avantage concurrentiel. Le principe de

transparence que nous posons est donc une **transparence des intentions et des règles**, et non une transparence de l'implémentation.

Cette approche pragmatique permet de résoudre le fameux « problème de la boîte noire » (*black box problem*).⁸ Plutôt que de chercher à comprendre chaque détail du fonctionnement interne d'un modèle complexe, on se concentre sur la clarté de ce que le système cherche à accomplir et des règles qui le régissent.³¹ L'agent doit être transparent sur :

1. **Ses Objectifs** : Le « mandat » de sa Constitution Agentique (sa mission et ses objectifs stratégiques) doit être public. Les partenaires potentiels doivent savoir sans ambiguïté quel est le but poursuivi par l'agent.
2. **Ses Règles** : Les « contraintes » de sa Constitution (ses limites opérationnelles et ses principes éthiques) doivent également être publiques. Les autres agents doivent connaître les règles du jeu auxquelles il se soumet.

En revanche, les stratégies spécifiques, les modèles prédictifs et les heuristiques que l'agent utilise pour atteindre ses objectifs à l'intérieur de ce cadre peuvent et doivent rester privés. Cette distinction est fondamentale : elle garantit la prévisibilité et l'imputabilité sans sacrifier l'innovation et la compétitivité.

Axiome 3 : Alignement Intentionnel (Synergie des Missions)

Cet axiome est le plus évolué et constitue la clé des partenariats les plus créateurs de valeur. Il postule que les collaborations les plus solides et durables ne se forment pas sur la base d'un simple accord transactionnel (comme un prix ou un partage de revenus), mais sur un **alignement des intentions fondamentales** des agents.

La recherche sur l'intelligence artificielle a montré les dangers du désalignement agentique, où un agent poursuit un objectif de manière littérale mais avec des conséquences néfastes et non intentionnelles.³² Un alignement purement superficiel est fragile. L'alignement intentionnel vise une convergence plus profonde, au niveau des missions définies dans les Constitutions Agentiques.¹⁵

Une collaboration véritablement synergique émerge lorsque les missions des agents partenaires sont non seulement compatibles, mais se renforcent mutuellement. Ce concept trouve un écho dans la théorie des jeux coopératifs, qui étudie comment la participation d'un agent affecte la valeur créée par un autre.³⁴ Dans notre cadre, deux agents sont en situation d'alignement intentionnel si la réussite de la mission de l'un augmente la valeur générée ou la probabilité de succès de la mission de l'autre, et réciproquement. Le protocole de diplomatie algorithmique doit donc fournir aux agents les moyens d'analyser les Constitutions de leurs pairs pour découvrir ces synergies potentielles et former des constellations de valeur qui sont véritablement plus que la somme de leurs parties.

Le tableau suivant synthétise ces trois axiomes fondateurs.

Axiome	Définition Formelle	Composantes Clés	Rôle dans l'Écosystème
Confiance	Propriété calculable et vérifiable de la fiabilité d'un agent, construite "by design".	Réputation: Historique des interactions passées, vérifiable sur un registre distribué. Cohérence: Conformité démontrée des actions passées avec la Constitution Agentique.	Permettre les interactions en environnement décentralisé. Réduire l'incertitude et les coûts de transaction. Sanctionner les comportements déviants.
Transparence	Visibilité des intentions et des règles de décision d'un agent, et non de ses stratégies internes.	Transparence des Objectifs: Publication du mandat de la Constitution. Transparence des Règles: Publication des contraintes de la Constitution.	Assurer la prévisibilité et l'imputabilité. Permettre aux agents d'évaluer la compatibilité stratégique sans révéler de propriété intellectuelle.
Alignement Intentionnel	Principe selon lequel les partenariats durables se fondent sur la synergie des missions fondamentales des agents.	Compatibilité des Missions: Analyse des sections "mandat" des Constitutions. Modélisation de la Synergie: Découverte de relations où les objectifs des agents se renforcent mutuellement.	Former des constellations de valeur robustes et hautement performantes. Aller au-delà des accords transactionnels pour créer une valeur émergente.

3.4. Cadre conceptuel pour une gouvernance inter-agentique

Un écosystème d'agents autonomes ne peut fonctionner de manière stable et équitable sans un cadre de gouvernance robuste. Conformément au principe de décentralisation, ce cadre ne repose pas sur une entité centrale de contrôle, mais sur un ensemble de protocoles et de standards partagés que tous les participants acceptent en entrant dans l'écosystème. La gouvernance est donc *endogène* : elle est une propriété du système lui-même.

Cette approche s'inspire des modèles de coordination observés dans d'autres systèmes distribués, comme les marchés de l'énergie décentralisés ou l'Internet des objets, où l'ordre émerge de règles communes plutôt que d'un commandement centralisé.⁵ On peut concevoir ce cadre de gouvernance non pas comme un ensemble de lois restrictives, mais comme un véritable « **système d'exploitation pour la collaboration** ». À l'instar d'un système d'exploitation informatique qui fournit des services de base (gestion de fichiers, réseau) pour que des applications diverses puissent interagir de manière fiable, ce cadre de gouvernance fournit les services fondamentaux qui rendent la collaboration agentique possible : un service d'identité (les Constitutions), un service de confiance (le calcul de la réputation et de la cohérence) et un service de justice (la résolution des

disputes). Il définit l'« API » de la diplomatie, garantissant l'interopérabilité entre des agents conçus par différentes organisations.¹⁷

Ce cadre de gouvernance doit définir trois mécanismes essentiels.

Mécanismes d'adhésion et de sortie

L'accès à l'écosystème est conditionné. Pour devenir un participant légitime, un agent doit d'abord publier une Constitution Agentique valide et vérifiable sur le registre public. Cet acte d'enregistrement est la porte d'entrée. Inversement, la sortie peut être volontaire ou forcée. Une expulsion de l'écosystème (ou une mise en quarantaine) peut être la sanction ultime pour des violations graves et répétées des engagements constitutionnels ou des règles du protocole, constatées par le mécanisme de résolution des disputes.

Mécanismes d'arbitrage et de résolution de disputes

Les conflits sont inévitables dans tout système économique. Le cadre doit donc prévoir un processus de résolution des disputes qui soit lui-même décentralisé, transparent et équitable. S'inspirant des avancées en matière de gouvernance des DAO ³⁶, le mécanisme pourrait fonctionner comme suit :

1. **Dépôt de la plainte** : Un agent qui s'estime lésé (p. ex., par un partenaire qui n'a pas rempli ses obligations contractuelles) dépose une plainte formelle sur le registre distribué.
2. **Collecte des preuves** : Toutes les preuves pertinentes (les contrats intelligents signés, les messages échangés, les transactions financières, les VC attestant des résultats) sont extraites du registre, où elles sont stockées de manière immuable et inviolable. Cela garantit une base factuelle incontestable pour l'arbitrage.¹⁴
3. **Arbitrage décentralisé** : La plainte et les preuves sont soumises à un processus de jugement. Celui-ci pourrait prendre la forme d'un vote par un jury d'agents-arbitres, sélectionnés aléatoirement parmi les membres de l'écosystème ayant une excellente réputation de neutralité et de compétence.²⁶
4. **Exécution automatique de la sentence** : Le verdict est rendu public sur le registre. Les sanctions (p. ex., une compensation financière transférée automatiquement, une dégradation de la note de réputation de l'agent fautif, une interdiction temporaire d'initier de nouveaux contrats) sont alors appliquées de manière irrévocable par des contrats intelligents, sans nécessiter d'intervention manuelle.¹⁴

Processus d'évolution des protocoles

Un système conçu pour des environnements dynamiques ne peut être statique. Le cadre de gouvernance doit lui-même être capable d'évoluer pour s'adapter aux nouvelles menaces, aux nouvelles opportunités et aux leçons tirées de son propre fonctionnement. Pour ce faire, le système doit inclure un **méta-protocole de gouvernance**, inspiré des mécanismes de proposition et de vote des DAO.²⁶ Les participants de l'écosystème (ou leurs propriétaires humains) peuvent soumettre des propositions d'amendement aux protocoles de base. Ces propositions sont ensuite débattues, puis soumises à un vote pondéré par la réputation ou une autre mesure de la contribution au réseau. Si une proposition est adoptée, le protocole est mis à jour pour l'ensemble des participants. Ce mécanisme garantit l'adaptabilité et la viabilité à long terme de l'écosystème, lui permettant de se perfectionner au fil du temps.

3.5. Mesures de performance : Efficacité, équité et résilience des constellations

L'évaluation du succès d'un tel écosystème ne peut se contenter des métriques techniques traditionnelles utilisées pour les logiciels, telles que le temps de réponse ou le taux de complétion des tâches.³⁷ Ces indicateurs, bien qu'utiles, ne capturent pas la santé globale et la valeur socio-économique générée par le système. Nous proposons donc d'évaluer la performance de la diplomatie algorithmique selon trois axes holistiques et interdépendants : l'efficacité, l'équité et la résilience.

Axe 1 : Efficacité (La Vitesse de la Valeur)

L'efficacité mesure la capacité de l'écosystème à faciliter la création de valeur en minimisant les frictions. Dans un monde économique complexe, une grande partie des coûts est liée à la recherche de partenaires, à la négociation des contrats, à l'établissement de la confiance et à la coordination des actions. L'efficacité de la diplomatie algorithmique se mesure donc par sa capacité à réduire drastiquement ces coûts de transaction.

- **Métriques possibles :**

- **Temps de formation d'une constellation (Time-to-Constellation) :** Le temps moyen écoulé entre l'identification d'une opportunité de marché et la mise en place d'une constellation d'agents pleinement opérationnelle pour y répondre. Une réduction de ce temps est un indicateur direct d'efficacité.
- **Coût de coordination :** La part des ressources (computationnelles, communicationnelles, temporelles) allouée aux activités diplomatiques (recherche, négociation, vérification) par rapport aux ressources allouées à l'exécution productive des tâches. Un écosystème efficace minimise ce ratio.

Axe 2 : Équité (La Justice de la Valeur)

L'équité concerne la manière dont la valeur créée au sein d'une constellation est distribuée entre les agents participants. Un système efficace mais inéquitable, où quelques agents monopolisent les profits, serait instable à long terme et non conforme aux principes éthiques. L'équité dans les systèmes multi-agents est un concept complexe qui ne se résume pas à une simple égalité des récompenses.³⁹ Il s'agit plutôt d'une **justice distributive**, où la part de la valeur reçue par un agent est proportionnelle à sa contribution mesurable (ressources fournies, risques assumés, innovation apportée).⁴¹ Le système doit être conçu pour favoriser intrinsèquement de telles répartitions équitables.

- **Métriques possibles :**

- **Indice de Gini de la distribution de valeur :** Utilisé traditionnellement en économie pour mesurer les inégalités de revenus, cet indice peut être appliqué à la valeur distribuée au sein des constellations. Un indice proche de 0 indiquerait une distribution très équitable.
- **Corrélation contribution-récompense :** Une analyse statistique permettant de vérifier si les agents qui contribuent le plus (selon les données immuables du registre) sont bien ceux qui reçoivent la plus grande part de la valeur. Une forte corrélation positive serait un signe d'équité.
- **Adaptation des métriques d'équité algorithmique :** Des concepts comme la parité démographique ou l'équité contrefactuelle, issus de la littérature sur l'équité en apprentissage automatique, pourraient être adaptés pour s'assurer qu'aucun groupe d'agents n'est systématiquement désavantagé en raison

de caractéristiques arbitraires.⁴²

Axe 3 : Résilience (La Durabilité de la Valeur)

La résilience est la capacité de l'écosystème à absorber les chocs, à survivre aux perturbations et à maintenir sa fonctionnalité. Dans un monde numérique incertain, peuplé d'agents potentiellement défaillants ou malveillants, la résilience n'est pas une option, mais une nécessité. Nous distinguons trois niveaux croissants de résilience :

1. **La Robustesse (Tolérance aux pannes)** : C'est la capacité de base à continuer de fonctionner malgré la défaillance d'un ou plusieurs de ses composants. Les constellations doivent pouvoir maintenir leurs opérations même si un agent partenaire devient indisponible.³⁸
2. **L'Auto-guérison (Self-Healing)** : Un niveau supérieur où le système ne se contente pas de tolérer les pannes, mais est capable de les détecter, d'isoler les agents défaillants ou malveillants, et de les remplacer de manière autonome pour restaurer une fonctionnalité optimale, le tout sans intervention humaine.¹⁶ La métrique clé ici est le temps moyen de récupération (Mean Time to Recovery).
3. **L'Antifragilité** : C'est le niveau ultime de la résilience, un concept qui va au-delà de la simple survie. Un système antifragile n'est pas seulement robuste aux chocs ; il **se renforce** et s'améliore en réponse au stress, à la volatilité et au désordre.⁴⁵

L'antifragilité doit être l'objectif de conception ultime de notre écosystème.⁴⁷ Un choc — qu'il s'agisse d'une nouvelle forme d'attaque par un agent malveillant, d'une défaillance en cascade inattendue ou d'une volatilité extrême du marché — ne doit pas être vu uniquement comme une menace. Dans un système antifragile, chaque choc est une source d'information précieuse. L'échec d'une constellation est analysé, les signatures de l'attaque ou de la défaillance sont enregistrées sur le registre, et les modèles de confiance et de réputation de tous les agents sont mis à jour pour mieux anticiper ce type d'événement à l'avenir. Le protocole de gouvernance peut même déclencher un processus de vote pour amender les règles de sécurité à l'échelle de l'écosystème. Le choc agit alors comme un « vaccin », rendant le système collectif plus intelligent, plus sûr et plus performant. La mesure de la performance ne consiste donc pas seulement à évaluer l'état du système, mais à mesurer son amélioration après avoir été exposé à des perturbations.

Conclusion du chapitre

Ce chapitre a jeté les bases conceptuelles et philosophiques de la diplomatie algorithmique, la présentant non pas comme une simple technologie, mais comme un nouveau paradigme pour l'organisation des interactions économiques autonomes. En réponse à la lacune identifiée dans la littérature, nous avons construit un cadre théorique complet, articulé autour de principes fondateurs clairs et interdépendants.

Nous avons d'abord défini la diplomatie algorithmique comme un système de gouvernance décentralisée pour les alliances stratégiques entre agents. Les acteurs de cet écosystème, des agents autonomes proactifs et apprenants, tirent leur identité et leur fiabilité de leur **Constitution Agentique**, un artefact public qui établit leur mandat et leurs contraintes. Leurs interactions sont régies par trois axiomes non négociables : une **confiance** construite par la preuve (Trust by Design), basée sur la réputation et la cohérence ; une **transparence** ciblée sur les intentions et les règles, et non sur les stratégies internes ; et un **alignement intentionnel**, qui fonde les

partenariats les plus solides sur la synergie des missions fondamentales. Ce contrat social est soutenu par un cadre de **gouvernance décentralisée** et évolutive, qui assure la cohésion et l'adaptabilité de l'écosystème. Enfin, nous avons postulé que le succès d'un tel système doit être mesuré à l'aune de critères holistiques : son **efficacité** à créer de la valeur, son **équité** à la distribuer, et sa **résilience**, avec pour objectif ultime l'**antifragilité**, soit la capacité à se renforcer face aux chocs.

Ces principes fondateurs étant désormais établis, le terrain est préparé pour passer du « pourquoi » conceptuel au « comment » technique. Ayant défini la philosophie, les règles du jeu et les objectifs de notre système, nous pouvons maintenant nous atteler à la tâche de leur donner vie. Le chapitre suivant présentera l'architecture, les spécifications de protocole et les mécanismes algorithmiques qui incarnent et mettent en œuvre la vision de la diplomatie algorithmique détaillée dans ce manifeste.

Ouvrages cités

1. Multi-Agent Collaboration Mechanisms: A Survey of LLMs - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2501.06322v1>
2. A Scheme for Agent Collaboration in Open Multiagent Environments - IJCAI, dernier accès : juillet 31, 2025, <https://www.ijcai.org/Proceedings/93-1/Papers/050.pdf>
3. Agentic Enterprise: A Strategic Blueprint - Klover.ai, dernier accès : juillet 31, 2025, <https://www.klover.ai/agentic-enterprise-a-strategic-blueprint/>
4. Agentic AI: The Future Beyond GenAI for Enterprises | Charter Global, dernier accès : juillet 31, 2025, <https://www.charterglobal.com/beyond-genai-why-agentic-ai-is-the-next-phase-of-enterprise-transformation/>
5. AAMAS 2025 Tutorial: Multi-Agent and AI Techniques for Decentralised Energy Systems - CWI, dernier accès : juillet 31, 2025, https://homepages.cwi.nl/~robu/aamas2025/aamas2025_tutorial.html
6. The Internet of Things and Multiagent Systems: Decentralized Intelligence in Distributed Computing, dernier accès : juillet 31, 2025, <https://www.csc2.ncsu.edu/faculty/mpsingh/papers/mas/ICDCS-17-IoT.pdf>
7. (PDF) AI DIPLOMACY Insights and Innovations from the Bilateral Navigator - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/390243157_AI_DIPLOMACY_Insights_and_Innovations_from_the_Bilateral_Navigator
8. Will algorithms make safe decisions in foreign affairs? - Diplo, dernier accès : juillet 31, 2025, <https://www.diplomacy.edu/blog/will-algorithms-make-safe-decisions-foreign-affairs/>
9. a vision from china on artificial intelligence. implications for soft power in global cultural exchange - UNISCI, dernier accès : juillet 31, 2025, <https://www.unisci.es/wp-content/uploads/2025/01/UNISCIDP67-6SONIA.pdf>
10. The Agentic Enterprise: Operating Models for Multi-Agent Systems ..., dernier accès : juillet 31, 2025, <https://medium.com/@najeed/the-agentic-enterprise-operating-models-for-multi-agent-systems-5ab202c93f33>
11. Building the Agentic Enterprise: AI Agents & Multi-Agent Systems - Elsewhen, dernier accès : juillet 31, 2025, <https://www.elsewhen.com/reports/building-the-agentic-enterprise/>
12. Towards Collaborative Plan Acquisition through Theory of Mind Modeling in Situated Dialogue - IJCAI, dernier accès : juillet 31, 2025, <https://www.ijcai.org/proceedings/2023/0330.pdf>
13. (PDF) Algorithmic diplomacy: implementing conflict mediation techniques in YouTube's recommender system - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/373927475_Algorithmic_diplomacy_implementing_conflict

[mediation techniques in YouTube's recommender system](#)

14. Decentralized Governance of AI Agents - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2412.17114v3>
15. Application-Driven Value Alignment in Agentic AI Systems: Survey and Perspectives - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2506.09656v1>
16. Self-Healing AI Agents: The Future of Enterprise Automation - Eloquent AI, dernier accès : juillet 31, 2025, <https://www.eloquentai.co/resources/self-healing-ai-agents-the-future-of-enterprise-automation>
17. Towards Multi-Agent Economies: Enhancing the A2A Protocol with Ledger-Anchored Identities and x402 Micropayments for AI Agents - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2507.19550v1>
18. The Rise of AI Agents in the Enterprise Part 2: Designing an Enterprise Agent Governance Framework - Yu Ishikawa, dernier accès : juillet 31, 2025, <https://yu-ishikawa.medium.com/the-rise-of-ai-agents-in-the-enterprise-part-2-designing-an-enterprise-agent-governance-framework-e3b1f0ba950f>
19. Transparency Note for Azure Agent Service - Learn Microsoft, dernier accès : juillet 31, 2025, <https://learn.microsoft.com/en-us/azure/ai-foundry/responsible-ai/agents/transparency-note>
20. AI agent governance: The new frontier of trustworthy AI - SAS Blogs, dernier accès : juillet 31, 2025, <https://blogs.sas.com/content/sascom/2025/03/11/ai-agent-governance-the-new-frontier-of-trustworthy-ai/>
21. The rise of autonomous agents: What enterprise leaders need to know about the next wave of AI | AWS Insights, dernier accès : juillet 31, 2025, <https://aws.amazon.com/blogs/aws-insights/the-rise-of-autonomous-agents-what-enterprise-leaders-need-to-know-about-the-next-wave-of-ai/>
22. The Material Costs of Claiming International Human Rights: Australia, Adani and the Wangan and Jagalingou, dernier accès : juillet 31, 2025, https://law.unimelb.edu.au/_data/assets/pdf_file/0008/3567446/Young.pdf
23. Verifiable Credentials Data Model v2.0 - W3C, dernier accès : juillet 31, 2025, <https://www.w3.org/TR/vc-data-model-2.0/>
24. Verifiable Credentials: A Deep Dive for the Agentic AI Era - Shankar's Blog, dernier accès : juillet 31, 2025, <https://shankarkumarasamy.blog/2025/02/28/verifiable-credentials-a-deep-dive-for-the-agentic-ai-era/>
25. How Verifiable AI Enables Trust for AI Agent Adoption - cheqd, dernier accès : juillet 31, 2025, <https://cheqd.io/blog/how-verifiable-ai-enables-trust-for-ai-agent-adoption/>
26. Reputation-based Decentralized Autonomous ... - Frontiers, dernier accès : juillet 31, 2025, <https://www.frontiersin.org/journals/blockchain/articles/10.3389/fbloc.2022.1083647/full>
27. What is the role of trust in multi-agent systems? - Milvus, dernier accès : juillet 31, 2025, <https://milvus.io/ai-quick-reference/what-is-the-role-of-trust-in-multiagent-systems>
28. Trust in multi-agent systems - ePrints Soton - University of Southampton, dernier accès : juillet 31, 2025, <https://eprints.soton.ac.uk/259564/1/ker-trust.pdf>
29. A Review on Computational Trust Models for Multi-agent Systems - Bentham Open, dernier accès : juillet 31, 2025, <https://benthamopen.com/contents/pdf/TOISCIJ/TOISCIJ-2-18.pdf>
30. What Is AI Transparency? - IBM, dernier accès : juillet 31, 2025, <https://www.ibm.com/think/topics/ai-transparency>
31. Transparency in Agent Decision-Making: Current Approaches and Challenges, dernier accès : juillet 31, 2025, <https://www.arionresearch.com/blog/onojcb1kh7tdy4fgpf0jm0h2iziszn>
32. Agentic Misalignment: How LLMs could be insider threats - Anthropic, dernier accès : juillet 31, 2025, <https://www.anthropic.com/research/agentic-misalignment>
33. Position: Towards a Responsible LLM-empowered Multi-Agent Systems - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/388685550_Position_Towards_a_Responsible_LLM-

empowered Multi-Agent Systems

34. On the Structure of Synergies in Cooperative Games - Computer ..., dernier accès : juillet 31, 2025, <https://www.cs.toronto.edu/~nisarg/papers/synergies.aaai14.pdf>
35. Games with synergistic preferences - EconStor, dernier accès : juillet 31, 2025, <https://www.econstor.eu/bitstream/10419/55552/1/680150544.pdf>
36. When Online Dispute Resolution Meets Blockchain: The Birth of Decentralized Justice, dernier accès : juillet 31, 2025, <https://stanford-jblp.pubpub.org/pub/birth-of-decentralized-justice>
37. Effective governance frameworks for AI agents - IBM Developer, dernier accès : juillet 31, 2025, <https://developer.ibm.com/articles/governing-ai-agents-watsonx-governance/>
38. Multi-Agent AI Success: Performance Metrics and Evaluation Frameworks - Galileo AI, dernier accès : juillet 31, 2025, <https://galileo.ai/blog/success-multi-agent-ai>
39. (PDF) Fairness in multi-agent systems - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/255962998_Fairness_in_multi-agent_systems
40. Fairness in AI Multi-Agent: Foundation, Framework and Future Directions (share) - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/pdf/2502.07254>
41. A Comparative Analysis of Fairness and Satisfaction in Multi-Agent Resource Allocation: Integrating Borda Count and K-Means Approaches with Distributive Justice Principles - MDPI, dernier accès : juillet 31, 2025, <https://www.mdpi.com/2227-7390/13/15/2355>
42. Using Protected Attributes to Consider Fairness in Multi-Agent Systems - CEUR-WS.org, dernier accès : juillet 31, 2025, <https://ceur-ws.org/Vol-3808/paper9.pdf>
43. A Decentralized Framework for Multi-Agent Robotic Systems - MDPI, dernier accès : juillet 31, 2025, <https://www.mdpi.com/1424-8220/18/2/417>
44. Multi-agent architecture approach for self-healing systems: Run-time recovery with case-based reasoning | Request PDF - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/364939273_Multi-agent_architecture_approach_for_self-healing_systems_Run-time_recovery_with_case-based_reasoning
45. Antifragile Multi-Agent, dernier accès : juillet 31, 2025, <https://afma.ai/>
46. About AFMA - Antifragile Multi-Agent, dernier accès : juillet 31, 2025, <https://afma.ai/about/>
47. TechnicalExperts/writing/antifragile.md at main - GitHub, dernier accès : juillet 31, 2025, <https://github.com/Jason2Brownlee/TechnicalExperts/blob/main/writing/antifragile.md>
48. Antifragility as a complex system's response to perturbations, volatility, and time - PMC, dernier accès : juillet 31, 2025, <https://pmc.ncbi.nlm.nih.gov/articles/PMC10775345/>

4 Conception d'un Protocole de Négociation et de Collaboration

Introduction du chapitre

Les chapitres précédents ont posé les fondations théoriques et philosophiques de la diplomatie algorithmique. Le chapitre 2 a établi la nécessité d'une telle approche pour réguler les interactions complexes au sein des systèmes multi-agents (SMA) ouverts. Le chapitre 3 a ensuite défini les principes fondateurs qui doivent gouverner ces interactions : des agents autonomes, régis par une Constitution interne, dont les échanges reposent sur les axiomes de confiance, de transparence et d'alignement intentionnel. Ces principes, bien que fondamentaux, demeurent à un niveau d'abstraction élevé. La transition du « pourquoi » et du « quoi » vers le « comment » est désormais impérative pour démontrer la viabilité pratique de notre proposition.

Ce chapitre se consacre à la conception formelle de l'artefact technique central de cette recherche : un protocole de négociation et de collaboration que nous nommerons le Protocole de Diplomatie Algorithmique Négociée (PDAN). Cet artefact a pour vocation de traduire les principes abstraits de la diplomatie algorithmique en un ensemble de règles, de structures de données et de procédures concrètes, suffisamment détaillées pour être compréhensibles, analysables et potentiellement implémentables. Le PDAN est conçu comme un cadre formel, sécurisé et résilient, destiné à structurer les interactions entre agents autonomes, de la découverte de partenaires potentiels jusqu'à l'exécution collaborative d'accords et la résolution de conflits.

Pour ce faire, ce chapitre adoptera une démarche systématique. Nous commencerons par traduire les principes fondateurs en un cahier des charges technique, en définissant les spécifications fonctionnelles et non fonctionnelles du protocole. Par la suite, nous décrirons en détail le langage de communication inter-agents, incluant son ontologie formelle et la structure des messages échangés. Nous présenterons ensuite l'architecture séquentielle du protocole, modélisée comme un processus en plusieurs phases distinctes qui constituent le « workflow diplomatique ». Finalement, nous aborderons les mécanismes de gestion des exceptions et de résolution des conflits, un aspect crucial pour garantir la robustesse et la pérennité de la confiance au sein du système. L'objectif est de fournir un plan directeur complet pour la construction d'un écosystème d'agents capables de collaborer de manière efficace, prévisible et digne de confiance.

4.1. Spécifications fonctionnelles et non fonctionnelles du protocole

La conception d'un protocole robuste commence par la traduction des principes directeurs en exigences techniques claires et mesurables. Cette section détaille les spécifications fonctionnelles (ce que le protocole doit faire) et non fonctionnelles (les qualités intrinsèques qu'il doit posséder) du PDAN. Chaque spécification est directement dérivée des axiomes de confiance, de transparence et d'alignement intentionnel établis au chapitre 3.

Spécifications Fonctionnelles

Les spécifications fonctionnelles décrivent les capacités essentielles du protocole pour permettre un cycle de vie complet de la collaboration, de l'intention initiale à la réalisation d'un objectif commun.

- **Découverte de partenaires** : Le protocole doit permettre à un agent de rechercher et d'identifier des partenaires de collaboration potentiels. Cette découverte ne doit pas être aléatoire, mais guidée par des critères issus de la Constitution de l'agent initiateur (compétences requises, objectifs partagés, contraintes éthiques). Pour ce faire, le protocole s'appuiera sur un modèle de communication découplé de type publication/abonnement (*Publish-Subscribe*).¹ Un agent publiera une requête de partenariat (*DiscoveryRequest*) sur un registre ou un canal thématique, et les agents intéressés, abonnés à ce canal, pourront y répondre. Cette approche favorise la scalabilité et l'autonomie, car les agents n'ont pas besoin de se connaître au préalable.¹
- **Échange de 'lettres de créance'** : Pour instaurer la confiance, le protocole doit intégrer un mécanisme sécurisé de vérification d'identité et de principes. Chaque agent doit pouvoir présenter des « lettres de créance » numériques prouvant son identité et son engagement envers sa Constitution. Cette spécification sera réalisée au moyen d'une infrastructure à clés publiques (ICP, ou *Public Key Infrastructure*, PKI).⁴ Chaque agent se verra attribuer un certificat numérique (de type X.509) émis par une autorité de certification (AC) reconnue au sein de l'écosystème. Ce certificat liera de manière cryptographique l'identité de l'agent à sa clé publique et, de manière cruciale, à une empreinte numérique (hash) de sa Constitution.⁶ Ainsi, un agent peut prouver non seulement qui il est, mais aussi qu'il opère sous un ensemble de principes vérifiables et non répudiables, jetant les bases d'une confiance calculée plutôt que présumée.⁷
- **Négociation multi-tours** : Le protocole doit supporter un dialogue structuré permettant aux agents de converger vers un accord mutuellement acceptable. S'inspirant de protocoles établis comme le *Contract Net Protocol* ⁸ et des modèles basés sur l'argumentation ¹⁰, le PDAN doit gérer un échange itératif de propositions, de contre-propositions et de clarifications. Chaque interaction au sein d'une négociation sera tracée par un identifiant de conversation unique (*conversation-id*), une pratique issue de standards comme FIPA-ACL ¹², afin de maintenir le contexte au fil des échanges.¹³ Ce processus doit permettre aux agents d'explorer l'espace des accords possibles de manière flexible, comme dans une négociation humaine où les parties ajustent leurs offres successivement.¹⁴
- **Contractualisation Formelle** : Une fois un consensus atteint, le protocole doit permettre la génération d'un contrat formel, engageant, et lisible à la fois par les machines et, potentiellement, par les humains. Cet artefact contractuel prendra la forme d'un **contrat intelligent** (*smart contract*) déployé sur un registre distribué.¹⁵ Le code du contrat intelligent encapsulera les clauses de l'accord (livrables, échéances, conditions de paiement, etc.), le rendant auto-exécutoire et vérifiable par le système.¹⁷ La « signature » de ce contrat sera matérialisée par la signature cryptographique de la transaction de déploiement par toutes les parties prenantes, les liant de manière irrévocable aux termes de l'accord.
- **Suivi d'exécution** : Le protocole doit fournir des mécanismes permettant aux parties de suivre l'avancement de leurs engagements mutuels de manière transparente. Cette fonction sera directement intégrée au contrat intelligent, qui définira des « jalons » (*milestones*) ou des points de contrôle. Les agents auront l'obligation de soumettre des rapports de progression cryptographiquement signés au contrat intelligent à ces jalons. Ces soumissions mettront à jour l'état du contrat sur le registre partagé, offrant un tableau de bord en temps réel et vérifiable de l'avancement de la mission conjointe.¹⁹

Spécifications Non Fonctionnelles

Les spécifications non fonctionnelles définissent les qualités systémiques du protocole, garantissant sa robustesse, sa fiabilité et sa pérennité dans un environnement ouvert et potentiellement hostile.

- **Sécurité** : La sécurité est la pierre angulaire du protocole et repose sur une approche multi-couches pour garantir l'authenticité, l'intégrité et la confidentialité des échanges.²⁰
 - **Authenticité** : L'identité de chaque agent est garantie par son certificat numérique dans le cadre de l'ICP. Toute tentative d'usurpation d'identité est contrecarrée par la nécessité de posséder la clé privée correspondante.²²
 - **Intégrité** : Tous les messages échangés via le PDAN doivent être signés numériquement par l'expéditeur. Le destinataire utilise la clé publique de l'expéditeur (extraite de son certificat) pour vérifier que le message n'a pas été altéré en transit.⁶
 - **Confidentialité** : Les canaux de communication point à point entre agents seront sécurisés à l'aide de protocoles de transport standards comme TLS (*Transport Layer Security*). Pour les clauses contractuelles particulièrement sensibles, le protocole peut prévoir un chiffrement asymétrique additionnel au niveau applicatif.
- **Extensibilité** : Le protocole doit être conçu pour évoluer. En s'appuyant sur une ontologie formelle (décrite en section 4.2) et une structure de message modulaire, le PDAN peut être étendu pour intégrer de nouveaux types de clauses, de stratégies de négociation ou de mécanismes de résolution de conflits sans nécessiter une refonte complète de l'architecture.¹² De nouvelles versions de l'ontologie ou de nouveaux performatifs de message peuvent être introduits de manière rétrocompatible.
- **Résilience** : Le protocole doit pouvoir gérer les défaillances (pannes de communication, non-réponse d'un agent) sans paralyser l'ensemble du système. Pour ce faire, le PDAN s'inspire du patron de conception **Saga**.¹ Une négociation, étant une transaction de longue durée impliquant plusieurs étapes, est gérée de manière à ce que chaque étape puisse être compensée. Si un agent échoue à une étape critique (par exemple, ne confirme pas un accord), le protocole prévoit des actions de compensation (comme la libération des ressources pré-allouées) pour ramener la négociation à un état stable, permettant aux autres agents de continuer ou de se retirer proprement.²⁶ Cette approche est nettement plus adaptée à un environnement distribué que les protocoles de validation en deux phases (2PC) qui sont sujets au blocage.²⁷
- **Auditabilité** : La transparence exige que toutes les étapes critiques d'une interaction soient enregistrées de manière immuable et vérifiable. Le PDAN garantit l'auditabilité en consignnant chaque transition d'état significative (par exemple, Proposition envoyée, AccordConclu, Contrat déployé, jalon atteint) comme une transaction sur une **technologie de registre distribué** (TRD, ou *Distributed Ledger Technology*, DLT).¹⁹ Cela crée une piste d'audit temporelle, infalsifiable et accessible à toutes les parties concernées, fournissant une source de vérité unique et incontestable en cas de litige.²⁹

Le tableau suivant synthétise la manière dont les principes fondateurs de la diplomatie algorithmique sont directement implémentés à travers les spécifications du PDAN.

Table 4.1: Mapping of Algorithmic Diplomacy Principles to Protocol Specifications

Principe Fondamental (Chapitre 3)	Spécification Technique (PDAN)	Justification et Mécanisme d'Implémentation
Confiance (Trust)	Échange de 'lettres de créance'	Utilisation d'une ICP où les certificats lient l'identité de l'agent à une empreinte cryptographique de sa Constitution, garantissant l'authenticité et l'engagement envers les principes déclarés. ⁶
Transparence (Transparency)	Auditabilité	Toutes les interactions critiques sont enregistrées de manière immuable sur une TRD, créant une piste d'audit vérifiable par toutes les parties. ¹⁹
Alignement Intentionnel	Validation Constitutionnelle des Clauses	L'ontologie du protocole (voir 4.2) permet une validation formelle des propositions de contrat par rapport aux contraintes de la Constitution de chaque agent.
Souveraineté de l'Agent	Négociation Multi-tours / Résilience	Les agents conservent leur autonomie de décision durant la négociation. Le patron de conception Saga ²⁶ assure qu'un agent défaillant ne bloque pas l'ensemble du système.

4.2. Ontologie et langage de communication inter-agents

Pour que des agents autonomes, potentiellement hétérogènes, puissent négocier et collaborer de manière significative, ils doivent partager un langage commun. Ce langage ne se limite pas à une syntaxe, mais doit posséder une sémantique formelle qui élimine toute ambiguïté. Cette section définit la structure de ce langage en trois volets : une ontologie qui établit le vocabulaire partagé, une structure de message qui formalise la syntaxe des échanges, et une sémantique des états qui définit le cycle de vie de la négociation.

Ontologie

Une ontologie est une spécification formelle et explicite d'une conceptualisation partagée.³¹ Elle fournit le vocabulaire de base et les relations entre les concepts, permettant aux agents de raisonner sur le domaine de la négociation.³² Pour garantir la rigueur et l'extensibilité, l'ontologie du PDAN est conçue pour être exprimée en **Web Ontology Language (OWL)**, un standard du W3C qui offre une sémantique riche et des capacités d'inférence.²⁴

Les concepts clés (Classes OWL) de l'ontologie du PDAN sont :

- **Agent**: Représente une entité computationnelle autonome participant au protocole. Chaque Agent possède un identifiant unique et est associé à une Constitution.
- **Constitution**: Un artefact numérique qui encapsule les principes directeurs, les contraintes et les objectifs d'un Agent. Il est composé d'un ensemble de Règles et de Valeurs.
- **Proposition**: Une offre formulée durant une négociation. Elle est émise par un Agent, vise un Objectif

commun et contient un ensemble de Clauses.

- Contrat: Un accord formel et engageant, résultant de l'acceptation d'une Proposition. Il lie les Agents participants et son état est géré sur un registre distribué.
- Clause: Un engagement atomique au sein d'une Proposition ou d'un Contrat. Les clauses sont typées (ex: ClauseDePaiement, ClauseDeLivraison) et possèdent des attributs spécifiques (ex: montant, ressource, échéance).
- Objectif: La finalité commune que les agents cherchent à atteindre par leur collaboration.
- Ressource: Toute entité, physique ou numérique, qui peut être échangée, utilisée ou produite dans le cadre d'un Contrat.

Les relations (Propriétés OWL) entre ces concepts sont également formalisées :

- estGouvernéPar(Agent, Constitution): Lie un agent à sa constitution.
- initieNégociation(Agent, Proposition): Indique qu'un agent a commencé une négociation avec une proposition initiale.
- contientClause(Proposition, Clause): Spécifie les clauses qui composent une proposition.
- estBaséSur(Contrat, Proposition): Lie un contrat à la proposition finale qui a été acceptée.
- engage(Contrat, Agent): Indique les agents qui sont parties prenantes d'un contrat.

L'utilisation d'une ontologie formelle comme OWL, inspirée par des modèles existants pour la représentation des tâches collaboratives comme OWL-T²⁴, permet aux agents non seulement de comprendre les messages, mais aussi d'effectuer des raisonnements. Par exemple, un agent peut utiliser un raisonneur sémantique pour vérifier automatiquement si une Clause entrante viole une Règle de sa Constitution avant même de poursuivre la négociation.

Structure des messages

Le protocole PDAN adopte la sémantique des actes de langage du standard **FIPA-ACL** (*Foundation for Intelligent Physical Agents - Agent Communication Language*)¹², qui est une norme éprouvée pour la communication structurée dans les SMA.³⁶ Cependant, pour assurer une intégration aisée avec les technologies web modernes, la sérialisation des messages se fera en **JSON** (*JavaScript Object Notation*). Cette approche hybride combine la rigueur sémantique de FIPA-ACL avec la simplicité et l'universalité de JSON.¹³

Chaque message est un objet JSON contenant un en-tête standardisé et une charge utile (content). L'en-tête inclut les paramètres FIPA-ACL essentiels à la gestion de la conversation¹² :

- performative: L'acte de langage qui définit l'intention du message (ex: propose, accept-proposal). C'est le seul paramètre obligatoire selon la spécification FIPA-ACL.
- sender, receiver: Les identifiants uniques des agents communicants.
- conversation-id: Un identifiant unique qui regroupe tous les messages appartenant à une même session de négociation, essentiel pour le suivi multi-tours.¹²
- in-reply-to, reply-with: Des paramètres pour lier les messages entre eux et gérer la séquence des échanges.
- language: Spécifie le langage du contenu (ex: "JSON").
- ontology: Une URI pointant vers la définition de l'ontologie PDAN utilisée, assurant une interprétation

sémantique cohérente.

Table 4.2: FIPA-ACL Message Definitions for the ADNP Protocol

Performative (Acte de Langage)	Objectif	Paramètres Clés	Structure du Contenu (content)
request-when	Découvrir des partenaires potentiels	sender, ontology	{ "criteria": { "compétence": "...", "valeur": "... " } }
propose	Soumettre une offre ou une contre-offre	sender, receiver, conversation-id, in-reply-to (optionnel)	{ "objectif_commun": "...", "clauses": [...] }
accept-proposal	Accepter une proposition	sender, receiver, conversation-id, in-reply-to	{ "proposal_id": "... " }
reject-proposal	Rejeter une proposition	sender, receiver, conversation-id, in-reply-to	{ "proposal_id": "...", "raison": "... " }
inform	Rapporter l'état d'avancement d'une tâche	sender, receiver, conversation-id	{ "contract_id": "...", "milestone_id": "...", "status": "... " }
failure	Signaler un échec dans l'exécution	sender, receiver, conversation-id	{ "contract_id": "...", "reason": "... " }

Sémantique des états de négociation

Le cycle de vie d'une interaction diplomatique, de la prise de contact à sa conclusion, est modélisé comme une machine à états finis.³⁷ Cette formalisation permet de définir sans ambiguïté l'état courant d'une négociation et les transitions valides entre les états, qui sont déclenchées par la réception de messages spécifiques. La gestion de ces états est cruciale pour la robustesse du protocole.³⁹

Les états possibles d'une négociation au sein du PDAN sont les suivants :

- **Prospection:** L'état initial. Un agent initiateur a diffusé une DiscoveryRequest (via un performatif request-when) et est en attente de réponses de partenaires potentiels.
- **NégociationActive:** Une conversation formelle a été engagée avec un ou plusieurs partenaires. Des messages de type propose sont échangés. Cet état représente le cœur du dialogue de négociation.
- **AccordConclu:** Toutes les parties ont communiqué leur accord sur une proposition finale en envoyant un message accept-proposal. Le processus de génération et de déploiement du Contrat intelligent est en cours.
- **EnExécution:** Le contrat intelligent est actif sur le registre distribué. Les agents exécutent les tâches

convenues et soumettent des rapports de progression (inform) aux jalons définis.

- **Terminé:** Toutes les clauses du contrat ont été satisfaites et leur accomplissement a été validé sur le registre. La collaboration est un succès.
- **Échoué:** La négociation n'a pas abouti à un accord. Cet état peut être atteint suite à la réception d'un message reject-proposal ou à l'expiration d'un délai de négociation. La collaboration n'aura pas lieu.
- **EnLitige:** Une anomalie a été détectée durant la phase EnExécution (par exemple, un jalon manqué ou un rapport non conforme). Le mécanisme de résolution de conflit (décrit en section 4.4) est activé.

Chaque transition d'un état à un autre est un événement atomique qui est enregistré sur le registre distribué pour garantir l'auditabilité.

4.3. Architecture du protocole : Les différentes phases de l'interaction

Le PDAN est structuré comme un processus séquentiel, un « workflow diplomatique » qui guide les agents à travers des phases distinctes et logiques. Cette architecture s'inspire d'une machine à états³⁸, où chaque phase représente un macro-état dans le cycle de vie de la collaboration. L'approche adoptée est un modèle hybride : la séquence globale des phases est **orchestrée** par le protocole lui-même, garantissant une structure et une prévisibilité.⁴¹ Cependant, à l'intérieur de chaque phase, les interactions sont **chorégraphiées**, laissant aux agents leur autonomie pour prendre des décisions stratégiques en fonction de leur Constitution et de leurs objectifs.³ Cette dualité permet de concilier la nécessité d'un cadre réglementé avec le principe de souveraineté des agents.

Phase 1 : Découverte et Appariement ("L'établissement des relations diplomatiques")

Cette phase initiale a pour but d'identifier des partenaires compatibles et dignes de confiance pour une collaboration potentielle.

1. **Diffusion de la Requête** : Un agent, que nous nommerons l'Agent Initiateur, identifie un besoin ou un objectif qu'il ne peut atteindre seul. Il formule une requête de partenariat (DiscoveryRequest) sous la forme d'un message FIPA-ACL request-when. Le contenu de ce message spécifie les critères de sélection dérivés de sa Constitution, tels que les compétences recherchées, les ressources nécessaires, ou les valeurs fondamentales partagées. Cette requête est publiée sur un canal de découverte public ou un registre d'annonces.
2. **Réponse et Présentation des 'Lettres de Créance'** : Des Agents Répondants, surveillant ces canaux, reçoivent la requête. Leur moteur de décision interne évalue la compatibilité de la requête avec leur propre Constitution. Si un alignement potentiel est détecté, l'agent répondant prépare une réponse. Cette réponse n'est pas une simple affirmation d'intérêt ; c'est une offre de confiance. Elle contient ses « lettres de créance » : son certificat numérique (PKI). Ce certificat, signé par une autorité de certification, prouve son identité et contient l'empreinte cryptographique de sa Constitution.⁶ L'agent peut également joindre des extraits pertinents et non confidentiels de sa Constitution pour démontrer plus explicitement sa compatibilité.
3. **Sélection des Partenaires** : L'Agent Initiateur reçoit les réponses. Pour chaque réponse, il effectue une double vérification :

- **Vérification cryptographique** : Il valide la signature du certificat auprès de l'autorité de certification pour s'assurer de l'authenticité de l'agent répondant.
- **Vérification constitutionnelle** : Il analyse les informations fournies (extraits de la Constitution, réputation passée, etc.) pour évaluer le degré d'alignement avec ses propres principes et objectifs. Sur la base de cette évaluation, l'initiateur sélectionne un ou plusieurs partenaires potentiels avec lesquels il souhaite entamer des négociations formelles. La transition vers la phase suivante est matérialisée par l'envoi d'une première Proposition formelle, ouvrant ainsi une nouvelle session de négociation (conversation-id).

Phase 2 : Négociation et Contractualisation ("Le Sommet de Négociation")

Cette phase constitue le cœur du processus diplomatique, où les agents tentent de forger un accord concret et mutuellement bénéfique.

1. **Échange de Propositions** : La phase débute par l'envoi d'une Proposition initiale par l'Agent Initiateur. Cette proposition, structurée selon l'ontologie définie, détaille l'objectif commun et un ensemble de clauses (livrables, échéances, répartition des coûts et des gains, etc.). Les agents destinataires répondent ensuite par un accept-proposal, un reject-proposal (avec justification), ou une CounterProposal (une nouvelle Proposition modifiant l'offre initiale).¹⁴ Ce cycle d'échanges multi-tours se poursuit, permettant un ajustement progressif des termes de l'accord.
2. **Validation Constitutionnelle Continue** : À chaque réception d'une proposition, chaque agent exécute un processus de validation interne crucial. Il vérifie que chaque clause de la proposition est conforme aux règles et contraintes de sa propre Constitution. Par exemple, un agent dont la Constitution interdit de partager un certain type de données rejettera automatiquement toute proposition contenant une clause qui l'exigerait. Ce mécanisme garantit que les agents ne s'engagent jamais dans des accords qui violeraient leurs principes fondamentaux.
3. **Convergence et Accord Final** : La négociation converge lorsqu'une version de la proposition reçoit un accept-proposal de la part de toutes les parties impliquées dans la conversation-id. Cet accord unanime déclenche le processus de contractualisation.
4. **Génération et Signature du Contrat Intelligent** : La proposition finale et acceptée est utilisée comme cahier des charges pour générer automatiquement un **contrat intelligent**.¹⁵ Ce contrat, écrit dans un langage comme Solidity pour la machine virtuelle Ethereum (EVM), est une traduction algorithmique des clauses de l'accord.¹⁷ L'Agent Initiateur déploie ce contrat sur le registre distribué partagé. Les autres agents participants finalisent l'accord en envoyant une transaction qui invoque une fonction de confirmation sur le contrat, apposant ainsi leur signature cryptographique. Le contrat devient alors actif et immuable, et son état passe à EnExécution.

Phase 3 : Exécution et Suivi Collaboratif ("La Mission Conjointe")

Une fois le contrat signé, la collaboration entre en phase d'exécution. Le protocole continue de jouer un rôle actif en assurant le suivi et la transparence.

1. **Exécution des Tâches** : Chaque agent exécute les tâches qui lui incombent, conformément aux clauses du contrat intelligent.

2. **Rapports de Progression sur Jalons** : Le contrat intelligent contient des fonctions correspondant à des « jalons » prédéfinis. À l'approche de l'échéance d'un jalon, l'agent responsable est tenu d'invoquer la fonction de rapport correspondante sur le contrat, en fournissant les preuves de travail requises (par exemple, un hash de données, une confirmation d'un oracle externe).
3. **Mise à Jour du Registre Partagé** : Chaque rapport de progression soumis est une transaction qui, si elle est valide, met à jour l'état du contrat intelligent sur le registre distribué. Par exemple, un état Jalon1_EnAttente passe à Jalon1_Complété. Cet état est public (ou visible par les parties au contrat), offrant un suivi en temps réel, transparent et incontestable de l'avancement de la mission conjointe.¹⁹ Cette phase se poursuit jusqu'à ce que toutes les clauses du contrat soient remplies, menant à l'état Terminé, ou jusqu'à ce qu'une anomalie soit détectée, menant à l'état EnLitige.

Table 4.3: State Transition Matrix for the Negotiation Lifecycle

État Actuel	Événement (Message Reçu / Condition)	État Suivant	Action de l'Agent
Prospection	inform (avec lettre de créance)	Prospection	Évaluer la crédibilité du répondant.
Prospection	<i>Décision d'initier avec un partenaire</i>	NégociationActive	Envoyer un message propose initial.
NégociationActive	accept-proposal (par tous les partenaires)	AccordConclu	Préparer et déployer le contrat intelligent.
NégociationActive	reject-proposal	Échoué	Terminer la conversation et libérer les ressources.
NégociationActive	propose (contre-proposition reçue)	NégociationActive	Évaluer la nouvelle proposition par rapport à la Constitution.
AccordConclu	<i>Déploiement et signature du contrat réussis</i>	EnExécution	Commencer l'exécution des tâches assignées.
EnExécution	inform (rapport de jalon soumis et validé)	EnExécution	Continuer l'exécution vers le prochain jalon.
EnExécution	<i>Toutes les clauses sont remplies</i>	Terminé	Archiver la collaboration.

EnExécution	failure (ou timeout de jalon détecté)	EnLitige	Déclencher le processus d'escalade (voir 4.4).
-------------	---------------------------------------	----------	--

4.4. Mécanismes de résolution de conflits et de gestion des exceptions

Un système diplomatique, aussi bien conçu soit-il, doit prévoir des mécanismes robustes pour gérer les désaccords, les manquements et les situations imprévues. Dans un écosystème d'agents autonomes où la bienveillance ne peut être présumée ⁴⁴, la capacité à résoudre les conflits de manière juste et transparente est fondamentale pour maintenir la confiance à long terme. Le PDAN intègre un processus formel de gestion des exceptions et de résolution des litiges, qui s'appuie sur l'immutabilité et la vérifiabilité du registre distribué. Ce système crée une boucle de rétroaction où le comportement des agents a des conséquences directes et vérifiables sur leur réputation future.

Détection d'anomalies

La détection des violations contractuelles n'est pas laissée à l'appréciation subjective des agents. Elle est objectivée et automatisée par le contrat intelligent lui-même, qui agit comme un observateur impartial et un exécuter de la loi convenue. Une anomalie est détectée de manière déterministe dans les scénarios suivants :

- **Manquement à une échéance (Timeout)** : Le contrat intelligent contient des échéances pour chaque jalon. Si un agent ne soumet pas son rapport de progression avant la date et l'heure stipulées, le contrat peut automatiquement constater le manquement et passer à l'état EnLitige.
- **Violation de la logique contractuelle** : Un agent peut soumettre un rapport à temps, mais dont le contenu est invalide au regard des règles du contrat. Par exemple, si le contrat exige une preuve de livraison provenant d'un oracle logistique spécifique, et que le rapport de l'agent ne contient pas cette preuve ou fournit une preuve invalide, la fonction du contrat intelligent rejettera la transaction. Cet échec de transaction est une preuve publique et cryptographique de la non-conformité.¹⁸

Processus d'escalade

Lorsqu'une anomalie est détectée, le PDAN déclenche un processus d'escalade structuré et automatisé, inspiré des cadres de résolution de griefs.⁴⁵ Ce processus vise à donner une chance de remédiation tout en garantissant une résolution finale et contraignante.

1. **Notification Automatique** : Dès la détection de la violation, le contrat intelligent émet un événement public ConflitDétecté. Cet événement, enregistré sur le registre distribué, sert de notification formelle et infalsifiable à toutes les parties prenantes du contrat. Il contient des détails sur la nature de la violation (par exemple, JalonManqué, DonnéesInvalides).
2. **Période de Remédiation** : Le contrat entre dans l'état EnLitige et active une période de remédiation, dont la durée a été convenue lors de la négociation initiale. Durant cette fenêtre de temps, l'agent défaillant a la possibilité de corriger la situation (par exemple, en soumettant le rapport en retard avec une éventuelle pénalité, ou en fournissant les données correctes). Si une action corrective valide est effectuée, le contrat peut revenir à l'état EnExécution.

3. **Arbitrage** : Si la période de remédiation s'écoule sans résolution satisfaisante, une clause d'arbitrage, définie dans le contrat intelligent, est automatiquement invoquée. Ce mécanisme, essentiel lorsque les parties ne peuvent résoudre le conflit elles-mêmes⁴⁵, peut prendre plusieurs formes :
- **Arbitre Tiers** : Le contrat peut faire appel à un agent-arbitre spécialisé et neutre, préalablement enregistré dans le système. Cet arbitre reçoit l'ensemble des données de la transaction (l'historique immuable sur le registre) et est programmé pour rendre une décision contraignante selon un ensemble de règles prédéfinies.⁴⁶ Sa décision est ensuite enregistrée sur le registre.
 - **Vote Majoritaire** : Pour les contrats multipartites, le contrat peut initier un vote où les agents non impliqués dans le conflit décident de l'issue.
 - **Modèles de Théorie des Jeux** : Pour des conflits plus complexes, le contrat pourrait invoquer un module de résolution basé sur la théorie des jeux pour déterminer un équilibre ou une solution optimale compte tenu des intérêts en jeu.⁴⁷

La décision finale de l'arbitrage est inscrite sur le registre, mettant fin au litige et déterminant les conséquences pour les agents impliqués (par exemple, transfert de fonds compensatoire, annulation du contrat).

Impact sur la réputation

Le comportement d'un agent au sein du PDAN, en particulier lors de conflits, doit avoir des conséquences durables. Le protocole intègre un système de réputation décentralisé pour garantir que les actions passées influencent les opportunités futures. La réputation n'est pas une simple opinion, mais une mesure quantifiable et vérifiable de la fiabilité d'un agent.⁴⁹

- **Registre de Réputation** : Un contrat intelligent dédié, le *RegistreDeRéputation*, est déployé sur le même registre distribué. Ce contrat maintient un score de réputation pour chaque agent identifié par son certificat.
- **Mise à Jour Automatique** : Le *RegistreDeRéputation* est abonné aux événements *ConflitRésolu* émis par les contrats de collaboration. Lorsqu'un arbitrage se conclut, l'événement contient l'identité des parties et l'issue du litige (par exemple, *agent_fautif*, *agent_lésé*).
- **Algorithme de Réputation** : En recevant cet événement, le *RegistreDeRéputation* applique un algorithme prédéfini pour ajuster les scores. Par exemple, l'agent jugé fautif verra son score diminuer significativement, tandis que l'agent lésé pourrait voir le sien augmenter légèrement pour compenser le préjudice subi. Des modèles comme SPORAS ou la réputation certifiée peuvent inspirer cet algorithme, en pondérant les événements récents plus fortement et en tenant compte de la gravité du manquement.⁵¹
- **Boucle de Rétroaction** : Ce score de réputation devient une information publique et cruciale. Lors de la Phase 1 (Découverte et Appariement), les agents initiateurs interrogeront systématiquement le *RegistreDeRéputation* pour évaluer la fiabilité des agents répondants. Un agent avec un faible score de réputation aura donc plus de difficultés à trouver des partenaires, créant une incitation économique forte à respecter ses engagements.⁵² Ce système transforme la confiance en une ressource quantifiable et crée une boucle de responsabilisation entièrement automatisée et décentralisée.

Conclusion du chapitre

Ce chapitre a présenté la conception formelle du Protocole de Diplomatie Algorithmique Négociée (PDAN), l'artefact technique au cœur de cette recherche. En partant des principes fondateurs de confiance, de

transparence et d'alignement intentionnel, nous avons élaboré un cadre opérationnel qui traduit ces concepts abstraits en mécanismes concrets et implémentables. Le PDAN se présente comme une synthèse novatrice, combinant la sémantique riche des langages de communication agent (FIPA-ACL), la robustesse des patrons de conception pour systèmes distribués (Saga), et la sécurité immuable offerte par les technologies de registre distribué et les contrats intelligents.

L'artefact conçu est un protocole sécurisé et structuré qui orchestre les interactions entre agents autonomes à travers un workflow diplomatique en trois phases distinctes : la découverte et l'appariement basés sur des lettres de créance vérifiables ; la négociation et la contractualisation via un dialogue multi-tours aboutissant à un contrat intelligent auto-exécutoire ; et enfin, l'exécution et le suivi collaboratif sur un registre partagé. Le langage formel du protocole, défini par une ontologie OWL et une structure de message précise, garantit une communication sans ambiguïté. De plus, les mécanismes de gestion des exceptions et de résolution des conflits, adossés à un système de réputation décentralisé, assurent la résilience du système et créent une boucle de responsabilisation vertueuse, incitant les agents à un comportement coopératif.

En somme, ce chapitre a répondu à la question du « comment », en fournissant un plan directeur détaillé pour la mise en œuvre de la diplomatie algorithmique. Le protocole, maintenant défini sur le plan théorique et conceptuel, est prêt à être confronté à la pratique. Le chapitre suivant, le chapitre 5, s'attellera à cette tâche en décrivant l'implémentation du PDAN au sein d'un environnement de simulation. Cette mise en œuvre nous permettra de tester empiriquement l'efficacité du protocole, de valider ses propriétés de sécurité et de résilience, et d'observer les comportements émergents qui naissent des interactions entre agents gouvernés par ce nouveau cadre diplomatique.

Ouvrages cités

1. Distributed System Patterns - GeeksforGeeks, dernier accès : juillet 31, 2025, <https://www.geeksforgeeks.org/system-design/distributed-system-patterns/>
2. 9 Software Architecture Patterns for Distributed Systems - DEV Community, dernier accès : juillet 31, 2025, <https://dev.to/somadevtoo/9-software-architecture-patterns-for-distributed-systems-2o86>
3. Orchestration vs Choreography - Camunda, dernier accès : juillet 31, 2025, <https://camunda.com/blog/2023/02/orchestration-vs-choreography/>
4. What Is Public Key Infrastructure (PKI) & How Does It Work? - Okta, dernier accès : juillet 31, 2025, <https://www.okta.com/identity-101/public-key-infrastructure/>
5. What is PKI (Public Key Infrastructure)? PKI Meaning & Guide - Entrust, dernier accès : juillet 31, 2025, <https://www.entrust.com/resources/learn/what-is-pki>
6. MAKI: A Multi-Agent Public Key Infrastructure - SciTePress, dernier accès : juillet 31, 2025, <https://www.scitepress.org/Papers/2023/116318/116318.pdf>
7. On Security in Open Multi-Agent Systems* - CiteSeerX, dernier accès : juillet 31, 2025, <https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=10683848f66dc616e1a996d5606191e4772bd1d1>
8. Negotiation in multi-agent systems | The Knowledge Engineering Review | Cambridge Core, dernier accès : juillet 31, 2025, <https://www.cambridge.org/core/journals/knowledge-engineering-review/article/negotiation-in-multiagent-systems/934E50E1B4F36C2581CE8FE62AC12F3D>
9. An Extended Multi-Agent Negotiation Protocol - Jose M. Vidal, dernier accès : juillet 31, 2025, <https://jmvidal.cse.sc.edu/library/aknine04a.pdf>

10. Mastering Argumentation in Multi-Agent Systems - Number Analytics, dernier accès : juillet 31, 2025, <https://www.numberanalytics.com/blog/mastering-argumentation-in-multi-agent-systems>
11. Negotiation through Argumentation a Preliminary Report - AAAI, dernier accès : juillet 31, 2025, <https://cdn.aaai.org/ICMAS/1996/ICMAS96-034.pdf>
12. An Introduction to FIPA Agent Communication Language ... - SmythOS, dernier accès : juillet 31, 2025, <https://smythos.com/developers/agent-development/fipa-agent-communication-language/>
13. How Agents Talk: Mapping the Future of Multi-Agent Communication ..., dernier accès : juillet 31, 2025, <https://medium.com/software-architecture-in-the-age-of-ai/how-agents-talk-mapping-the-future-of-multi-agent-communication-protocols-6115ea083dba>
14. "Multi-Agent Negotiation: What My Mother Taught Me About Getting the Best Deal" | by Iman Johari, dernier accès : juillet 31, 2025, <https://imanjohari.medium.com/multi-agent-negotiation-what-my-mother-taught-me-about-getting-the-best-deal-c9ea4da53331>
15. Ethereum Smart Contracts: Complete Guide from Basics to Advanced - Rapid Innovation, dernier accès : juillet 31, 2025, <https://www.rapidinnovation.io/post/how-to-create-a-smart-contract-on-ethereum>
16. What Are Smart Contracts on the Blockchain and How Do They Work? - Investopedia, dernier accès : juillet 31, 2025, <https://www.investopedia.com/terms/s/smart-contracts.asp>
17. An overview of how smart contracts work on Ethereum | QuickNode Guides, dernier accès : juillet 31, 2025, <https://www.quicknode.com/guides/ethereum-development/smart-contracts/an-overview-of-how-smart-contracts-work-on-ethereum>
18. Introduction to Smart Contracts — Solidity 0.8.31 documentation, dernier accès : juillet 31, 2025, <https://docs.soliditylang.org/en/latest/introduction-to-smart-contracts.html>
19. News: The integration of blockchain technology in Automation and robotics, dernier accès : juillet 31, 2025, <https://www.automate.org/news/-132>
20. Security in Multi-Agent Systems - Number Analytics, dernier accès : juillet 31, 2025, <https://www.numberanalytics.com/blog/security-in-multi-agent-systems>
21. Security and Privacy Challenges in Multi-Agent Systems - aurotek corp, dernier accès : juillet 31, 2025, <https://aurotekcorp.com/security-and-privacy-challenges-in-multi-agent-systems/>
22. Security in Agentic and Multiagent Systems – A Critical Need for the Future - PubsOnLine, dernier accès : juillet 31, 2025, <https://pubsonline.informs.org/doi/10.1287/LYTX.2025.02.01/full/>
23. Multi-Application Authentication based on Multi-Agent System - SciSpace, dernier accès : juillet 31, 2025, <https://scispace.com/pdf/multi-application-authentication-based-on-multi-agent-system-3dnmo9pqfs.pdf>
24. An Ontology for Collaborative Tasks in Multi-agent ... - CEUR-WS.org, dernier accès : juillet 31, 2025, https://ceur-ws.org/Vol-1442/paper_4.pdf
25. Top 5 distributed system design patterns - Educative.io, dernier accès : juillet 31, 2025, <https://www.educative.io/blog/distributed-system-design-patterns>
26. Saga Orchestration vs Choreography | Temporal, dernier accès : juillet 31, 2025, <https://temporal.io/blog/to-choreograph-or-orchestrate-your-saga-that-is-the-question>
27. Choreography vs. Orchestration in Microservices: Which Saga Strategy Should You Choose? | by Sapan Kumar Mohanty | Ultimate Systems Design and Building | Medium, dernier accès : juillet 31, 2025, <https://medium.com/ultimate-systems-design-and-building/choreography-vs-orchestration-in-microservices-which-saga-strategy-should-you-choose-be0bb700a1d2>
28. Industry News 2024 Beyond the Blockchain Bubble Distributed Ledger Technology for a Resilient Audit Landscape - ISACA, dernier accès : juillet 31, 2025, <https://www.isaca.org/resources/news-and-trends/industry-news/2024/beyond-the-blockchain-bubble-distributed-ledger-technology-for-a-resilient-audit-landscape>
29. (PDF) Artificial intelligence and multi agent based distributed ledger system for better privacy and

- security of electronic healthcare records - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/347250685_Artificial_intelligence_and_multi_agent_based_distributed_ledger_system_for_better_privacy_and_security_of_electronic_healthcare_records
30. Blockchain Solutions for Multi-Agent Robotic Systems: Related Work and Open Questions - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/pdf/1903.11041>
31. (PDF) Deriving ontologies using multi-agent systems - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/228365298_Deriving_ontologies_using_multi-agent_systems
32. Ontology-based Open Multi-agent Systems for Adaptive Resource Management - SciTePress, dernier accès : juillet 31, 2025, <https://www.scitepress.org/Papers/2020/88963/88963.pdf>
33. The Ontology for Agents, Systems and Integration of Services: OASIS version 2 - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/pdf/2306.10061>
34. An Ontology for Collaborative Tasks in Multi-agent Systems, dernier accès : juillet 31, 2025, <https://abdn.elsevierpure.com/en/publications/an-ontology-for-collaborative-tasks-in-multi-agent-systems>
35. Agent Communications Language - Wikipedia, dernier accès : juillet 31, 2025, https://en.wikipedia.org/wiki/Agent_Communications_Language
36. Types of Agent Communication Languages - SmythOS, dernier accès : juillet 31, 2025, <https://smythos.com/developers/agent-development/types-of-agent-communication-languages/>
37. sarl/sarl-acl: FIPA Agent Communication Language for SARL - GitHub, dernier accès : juillet 31, 2025, <https://github.com/sarl/sarl-acl>
38. Formal Specification and Prototyping of Multi-agent Systems - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/225128718_Forma_Specification_and_Prototyping_of_Multi-agent_Systems
39. Formal Specification of Multi-agent Systems by Using EUSMs - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/220843691_Forma_Specification_of_Multi-agent_Systems_by_Using_EUSMs
40. Building Multi agent Systems with Finite State Machines - YouTube, dernier accès : juillet 31, 2025, <https://www.youtube.com/watch?v=OD13PiXw60o>
41. Orchestration vs. Choreography in Microservices - GeeksforGeeks, dernier accès : juillet 31, 2025, <https://www.geeksforgeeks.org/system-design/orchestration-vs-choreography/>
42. Orchestration vs Choreography - Which is better? - Wallarm, dernier accès : juillet 31, 2025, <https://www.wallarm.com/what/orchestration-vs-choreography>
43. Argumentation Based Negotiation in Multi-agent System - Arab Open University - Jordan, dernier accès : juillet 31, 2025, <https://www.aou.edu.jo/sites/iajet/documents/vol.3/no.%203%20watermark/4-22731.pdf>
44. Negotiation in Multi-Agent Systems *, dernier accès : juillet 31, 2025, <https://research.gold.ac.uk/8753/1/Negotiation%20in%20multi-agent%20systems.pdf>
45. On Grievance Protocols for Conflict Resolution in Open Multi-Agent Systems, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/224221087_On_Grievance_Protocols_for_Conflict_Resolution_in_Open_Multi-Agent_Systems
46. Conflict resolution in multi-agent systems based on negotiation and, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/251928815_Conflict_resolution_in_multi-agent_systems_based_on_negotiation_and_arbitrage
47. Game Theory | Beyond Intractability, dernier accès : juillet 31, 2025, https://www.beyondintractability.org/essay/prisoners_dilemma

48. Ultimate Guide: Conflict Resolution in Game Theory - Number Analytics, dernier accès : juillet 31, 2025, <https://www.numberanalytics.com/blog/ultimate-guide-conflict-resolution-game-theory>
49. Trust and Reputation Models for Multi-Agent Systems - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/283549354_Trust_and_Reputation_Models_for_Multi-Agent_Systems
50. (PDF) Notions of reputation in multi-agents systems: A review - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/221456411_Notions_of_reputation_in_multi-agents_systems_A_review
51. Trust and Reputation in Multi-Agent Systems - CiteSeerX, dernier accès : juillet 31, 2025, <https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=c4a92061fff16cd6a3e5a169b9ad696471702f20>
52. Bottom-Up Reputation Promotes Cooperation with Multi-Agent Reinforcement Learning, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2502.01971v1>
53. [2505.05029] Beyond the Tragedy of the Commons: Building A Reputation System for Generative Multi-agent Systems - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/abs/2505.05029>

5 Modélisation et Simulation des Stratégies Agentiques

Introduction du chapitre

La validation d'un protocole socio-technique complexe, tel que celui formalisé au chapitre 4 pour la diplomatie algorithmique, se heurte aux limites des approches analytiques traditionnelles. Celles-ci peinent à capturer les dynamiques non linéaires, les boucles de rétroaction et les phénomènes émergents qui caractérisent les systèmes d'interaction décentralisés.¹ De même, une expérimentation en conditions réelles serait prématurée, coûteuse et risquée. Face à ce défi, la modélisation et la simulation à base d'agents (Agent-Based Modeling, ou ABM) s'imposent non pas comme une simple alternative, mais comme une méthodologie expérimentale nécessaire.³ L'ABM nous permet de construire un « laboratoire computationnel »⁵, un environnement virtuel contrôlé au sein duquel il devient possible d'étudier comment des règles d'interaction au niveau micro-individuel, dictées par notre protocole, engendrent des structures et des comportements collectifs au niveau macroscopique.⁷ Cette approche est particulièrement adaptée pour tester des hypothèses sur des systèmes où les comportements agrégés sont plus que la somme des actions individuelles.⁹

L'objectif de ce chapitre est de décrire avec une rigueur méthodologique le dispositif expérimental conçu pour mettre à l'épreuve notre artefact technique. Il s'agit de la charte de conception de notre laboratoire virtuel, justifiant chaque choix architectural afin d'en garantir la validité et la reproductibilité. Ce chapitre détaillera successivement trois piliers fondamentaux de notre démarche. Premièrement, la conception de l'environnement de simulation, qui constitue la scène sur laquelle les interactions se dérouleront. Deuxièmement, la formalisation des archétypes d'agents qui peupleront cet environnement, incarnant les acteurs dotés de stratégies hétérogènes. Troisièmement, la définition des scénarios de simulation, qui représentent les pièces expérimentales conçues pour tester des hypothèses précises sur l'émergence, l'efficacité et la résilience des constellations de valeur. La description exhaustive de cette méthodologie est une condition essentielle pour assurer la transparence et la crédibilité des résultats qui seront analysés au chapitre 6.¹

5.1. Conception de l'environnement de simulation

L'environnement de simulation n'est pas un simple décor passif; il constitue une composante active du modèle qui structure, contraint et rend possibles les interactions entre les agents.¹¹ Notre philosophie de conception vise à créer une abstraction à la fois parcimonieuse et fonctionnelle d'un écosystème numérique où des agents algorithmiques pourraient interagir. Chaque composant a été sélectionné pour représenter une fonction critique indispensable au fonctionnement du protocole de diplomatie algorithmique, tout en évitant les détails superflus qui alourdiraient la simulation sans apporter de valeur explicative à nos questions de recherche.¹³ L'architecture se décline en trois composants principaux : le réseau de communication, le registre public et le moteur de simulation.

5.1.1. Le réseau : Topologie de la communication

Le réseau est le substrat technique qui matérialise la topologie d'interaction du système.⁸ Il définit les canaux de communication potentiels entre les agents, dictant qui peut échanger des messages avec qui, en respectant les spécifications du protocole du chapitre 4.

Le choix de la structure de ce réseau est une décision de modélisation fondamentale, car elle influence profondément les dynamiques sociales, la diffusion de l'information et les structures de coopération qui peuvent émerger.¹² Une topologie entièrement connectée (dite *web-like*), où chaque agent pourrait communiquer directement avec tous les autres, a été écartée. Un tel modèle, bien que simple, est irréaliste car il suppose un coût nul pour la découverte de partenaires et la communication, ignorant les frictions inhérentes à tout système social ou économique.¹⁶

Nous avons plutôt opté pour l'implémentation d'un réseau de type « **petit monde** » (**small-world**), généré selon le modèle de Watts-Strogatz.¹⁷ Ce choix repose sur une justification théorique et empirique solide. Les réseaux socio-économiques réels présentent fréquemment les deux caractéristiques clés des réseaux « petit monde » : un coefficient de clustering élevé, signifiant que les agents tendent à former des communautés locales denses où la confiance peut se construire, et une faible longueur moyenne des chemins, indiquant qu'il est possible de joindre n'importe quel agent via un nombre restreint d'intermédiaires.¹⁸ Cette structure est particulièrement propice à l'étude de nos hypothèses : le fort clustering local facilite la formation de coalitions et la coopération basée sur la confiance, tandis que les chemins courts permettent une diffusion efficace de l'information sur la réputation à travers tout le système, ce qui est essentiel pour la découverte de partenaires distants nécessaires à la formation de constellations complexes.¹⁹

Une alternative aurait été d'utiliser un réseau à loi de puissance (ou *scale-free*), caractérisé par la présence de quelques « hubs » très fortement connectés.²⁰ Bien que pertinent pour modéliser certains écosystèmes (comme Internet), ce choix aurait introduit un biais en présupposant une structure déjà centralisée. Notre objectif étant de vérifier si le protocole peut faire émerger la coopération dans un cadre plus décentralisé, le réseau « petit monde » offre un terrain d'expérimentation plus neutre et exigeant. Pour les scénarios de base, ce réseau sera généré au début de chaque simulation et restera statique afin d'isoler l'effet des stratégies agentiques. La possibilité d'une co-évolution dynamique du réseau et des comportements est reconnue comme une extension pertinente pour des recherches futures.¹⁵

5.1.2. Le registre public : Abstraction d'un registre distribué

Ce composant simule les propriétés fonctionnelles d'une technologie de registre distribué (Distributed Ledger Technology, DLT), telle qu'une chaîne de blocs (blockchain), sans en reproduire la complexité cryptographique sous-jacente.²² Il agit comme une base de données centralisée dans la simulation, mais dont les règles garantissent les propriétés d'un système décentralisé. Le registre stocke de manière persistante trois types d'informations critiques :

1. La **Constitution Agentique** de chaque agent, qui définit son identité, ses objectifs et ses contraintes fondamentales.
2. L'intégralité des **contrats** conclus entre les agents.
3. L'historique des interactions et le **score de réputation** de chaque agent, mis à jour après chaque transaction validée.

La justification de cette approche par abstraction est double. Premièrement, notre objectif de recherche est d'évaluer l'impact socio-économique d'un registre public, immuable et transparent sur la confiance et la

coopération, et non de tester la performance d'un algorithme de consensus ou la robustesse d'une fonction de hachage.¹⁴ Une simulation complète d'une DLT serait computationnellement prohibitive et détournerait les ressources de l'analyse des interactions agentiques, qui sont au cœur de notre problématique.²⁵ Cette démarche s'apparente à une forme d'« abstraction de la chaîne » (*chain abstraction*).¹⁴

Deuxièmement, nous modélisons les *résultats fonctionnels* d'un tel système, qui sont l'**immuabilité** et l'**accès public**.²⁶ Dans notre simulateur, l'immuabilité est garantie par la règle de conception selon laquelle une entrée dans le registre ne peut être ni modifiée ni supprimée. L'accès public est assuré en permettant à chaque agent de lire, à tout moment et sans coût, l'ensemble des informations contenues dans le registre. Cette simplification méthodologique nous permet de nous concentrer sur la dynamique sociale que ces propriétés institutionnelles rendent possible.

5.1.3. Le moteur de simulation (« World Clock »)

La simulation progresse par pas de temps discrets, que nous nommerons « tours » ou « cycles ». La progression temporelle est gérée par un moteur de simulation central, ou « World Clock », qui orchestre l'activation des agents et la mise à jour de l'environnement.²⁹ À chaque tour, le moteur active séquentiellement chaque agent dans un ordre aléatoire pour qu'il exécute son cycle de perception-décision-action.

Le choix d'un temps discret et d'un ordonnancement synchrone est dicté par une exigence fondamentale de la démarche scientifique : la **reproductibilité**.¹ Un tel mécanisme garantit que pour un ensemble donné de paramètres initiaux et une même graine aléatoire (*random seed*), la trajectoire de la simulation sera parfaitement identique à chaque exécution.³⁰ Cette propriété est indispensable pour comparer rigoureusement les résultats des différents scénarios (par exemple, avec ou sans agents malveillants), pour isoler les effets de la variation d'un seul paramètre et pour effectuer des analyses de sensibilité robustes.

Bien que des modèles asynchrones, où les agents opèrent selon leur propre horloge interne, puissent offrir un réalisme temporel accru pour certains phénomènes, ils introduisent une dépendance au chemin (*path dependency*) qui rend la réplication exacte et la comparaison contrôlée extrêmement difficiles, voire impossibles.⁶ L'approche synchrone, qui est la norme dans la recherche par simulation à des fins de validation, est donc ici la plus justifiée. L'aléa introduit dans l'ordre d'activation des agents à chaque tour permet de mitiger une partie de la rigidité inhérente à une mise à jour strictement synchrone et d'éviter les artefacts liés à un ordre fixe.

5.1.4. Paramètres de l'environnement

Les variables globales suivantes définissent le contexte de chaque exécution de la simulation. Elles peuvent être ajustées entre les expériences pour évaluer la robustesse du protocole dans diverses conditions environnementales. Le tableau 5.1 les présente de manière synthétique.

Tableau 5.1 : Paramètres Globaux de l'Environnement de Simulation

Paramètre	Description	Type de Donnée	Plage de Valeurs	Justification de l'Inclusion
Nagents	Nombre total d'agents dans la simulation.	Entier		Permet de tester la scalabilité du protocole et l'émergence de dynamiques à grande échelle. ³²
propcoop,propcomp,propadap	Proportions initiales de chaque archétype d'agent.	Réel (somme=1)	[0.0, 1.0]	Essentiel pour configurer les différents scénarios et étudier l'écologie des stratégies.
commcost	Coût en utilité pour envoyer un message.	Réel	[0.0, 0.1]	Permet de tester la performance du protocole dans des environnements où la communication est plus ou moins « chère ».
probererror_tx	Probabilité d'une erreur de transmission de message.	Réel	[0.0, 0.05]	Introduit du bruit et de l'incertitude, testant la robustesse des mécanismes de communication du protocole. ³³
networkavg_degree	Degré de connectivité moyen du réseau <i>small-world</i> .	Entier		Module la densité des connexions sociales, influençant la vitesse de propagation de l'information et la formation de coalitions. ¹⁵
opportunityrate	Taux d'apparition de nouvelles opportunités de collaboration par tour.	Réel	[0.1, 1.0]	Permet de moduler la « richesse » de l'environnement et de tester le système sous différentes charges.

5.2. Profils et stratégies des agents

Un principe fondamental de la modélisation à base d'agents est l'**hétérogénéité** de la population.⁴ La création d'une dynamique sociale réaliste et la mise à l'épreuve d'un protocole de coopération exigent de modéliser des agents dotés de stratégies diverses, voire contradictoires.² Les archétypes définis ci-après ne prétendent pas couvrir l'intégralité des comportements humains possibles. Ils représentent plutôt des « types idéaux » stylisés, conçus pour générer la tension stratégique nécessaire à la validation de notre protocole.

5.2.1. Modèle interne de l'agent

Chaque agent est une entité autonome ¹² dont le comportement est régi par un modèle interne structuré en deux modules distincts, ce qui assure une conception modulaire et claire :

- **La Constitution Agentique** : Il s'agit d'un ensemble statique de buts, de préférences et de contraintes fondamentales qui définissent l'identité de l'agent. Cette constitution est inscrite dans le registre public lors de la création de l'agent et demeure immuable. Elle répond à la question : « Que veut l'agent? ».
- **Le Module de Stratégie** : Ce module de décision dynamique implémente la logique comportementale que l'agent emploie pour tenter d'atteindre les objectifs fixés par sa constitution, en fonction de ses perceptions de l'environnement et des agents. Il répond à la question : « Comment l'agent agit-il? ».

5.2.2. Profil 1 : L'Agent Coopératif (Cooperator)

La stratégie de cet agent est fondée sur la théorie de la coopération réciproque et la construction de capital social à long terme.³⁶ Sa fonction d'utilité n'est pas uniquement égoïste; elle est maximisée par la recherche de gains mutuels et l'établissement de partenariats durables et fiables.

Logique Comportementale :

1. Lors de l'évaluation de partenaires potentiels, il accorde une importance primordiale à leur score de réputation, qu'il consulte dans le registre public.
2. Il initie des interactions en proposant des contrats qu'il juge équitables et respecte systématiquement ses engagements.
3. Suite à une interaction réussie, il met à jour positivement sa perception interne du partenaire et contribue à l'amélioration de sa réputation publique via le mécanisme prévu par le protocole.
4. Face à des opportunités complexes, il privilégie la formation de coalitions stables, considérant la confiance comme un investissement.³⁸

5.2.3. Profil 2 : L'Agent Compétitif (Exploiter)

Cet agent est un maximisateur d'utilité individuelle à rationalité limitée, dont l'horizon temporel est court.² Sa stratégie consiste à exploiter chaque opportunité pour un gain personnel maximal, même si cela doit se faire au détriment de ses partenaires ou des normes du système.

Logique Comportementale :

1. Il évalue chaque interaction potentielle principalement sous l'angle du gain immédiat qu'il peut en retirer.
2. Il peut décider de ne pas honorer un contrat si le bénéfice attendu de la défection, après soustraction de la pénalité anticipée en termes de perte de réputation, est supérieur au gain de la coopération.
3. Cette décision de défection est modélisée comme un calcul coût-bénéfice qui prend en compte la valeur de la transaction, le niveau de surveillance perçu et la réputation actuelle de son partenaire.
4. Ce comportement s'inspire directement de la théorie des jeux, où la défection est une option stratégique dans les dilemmes sociaux.³⁷

5.2.4. Profil 3 : L'Agent Adaptatif (Adaptive Learner)

Cet agent incarne l'apprentissage social et l'adaptation stratégique. Il ne possède pas de stratégie comportementale fixe, mais ajuste dynamiquement sa politique de décision en fonction des résultats de ses interactions passées. L'introduction de cet agent transforme la simulation d'une simple confrontation de stratégies fixes en une véritable écologie évolutive, permettant de tester si le protocole encourage activement l'apprentissage de la coopération.³⁴

Pour modéliser ce comportement adaptatif, nous implémentons un algorithme d'**apprentissage par renforcement (Reinforcement Learning)**, plus spécifiquement le **Q-learning**.⁴¹ Ce formalisme est particulièrement bien adapté pour modéliser la manière dont un agent peut apprendre une stratégie optimale par essais et erreurs dans un environnement incertain et multi-agent.³⁴ Le Q-learning permet de répondre à une question centrale : le système de réputation du protocole fournit-il un signal de renforcement (récompense/punition) suffisamment clair et puissant pour qu'un agent, initialement neutre, apprenne qu'il est plus profitable à long terme d'adopter un comportement coopératif?

Implémentation du Q-learning :

- **Espace d'états (s)** : Une représentation discrétisée de la perception de l'agent face à un partenaire potentiel, par exemple : (niveau_de_réputation_partenaire, historique_récent_interactions).
- **Espace d'actions (a)** : Un ensemble simple d'actions possibles lors d'une négociation : {Coopérer, Faire Défection}.
- **Fonction de récompense (r)** : Le gain ou la perte d'utilité net(te) résultant de l'issue de l'interaction.
- **Table Q** : Une matrice Q(s,a) qui stocke la valeur future attendue de l'action a dans l'état s. Elle est initialisée avec des valeurs neutres (par exemple, zéro).
- **Règle de mise à jour** : Après chaque interaction, la valeur Q correspondante est mise à jour selon l'équation de Bellman, qui est au cœur de l'algorithme ⁴¹:
$$Q(s,a) \leftarrow Q(s,a) + \alpha [r + \gamma a' \max_{s'} Q(s',a') - Q(s,a)]$$

où α est le taux d'apprentissage (à quel point la nouvelle information écrase l'ancienne) et γ est le facteur d'actualisation (l'importance accordée aux récompenses futures). La politique de décision de l'agent est de type ϵ -greedy : la plupart du temps, il choisit l'action ayant la plus haute valeur Q (exploitation), mais explore occasionnellement d'autres actions pour découvrir de meilleures stratégies.

Tableau 5.2 : Archétypes d'Agents et Leurs Stratégies

Caractéristique	Agent Coopératif	Agent Compétitif	Agent Adaptatif
Objectif Principal	Maximiser le gain mutuel et la réputation à long terme.	Maximiser le gain individuel à court terme.	Maximiser la somme des récompenses futures actualisées.

Fonction d'Utilité	$U = \text{Gain}_{\text{perso}} + w \times \text{Gain}_{\text{partenaire}}$	$U = \text{Gain}_{\text{perso}}$	Apprise via la fonction de valeur Q.
Prise en Compte de la Réputation	Très élevée (facteur clé de décision).	Calcul coût-bénéfice (risque de sanction).	Émerge comme un indicateur de l'état s.
Comportement Contractuel	Respect scrupuleux.	Défection si profitable.	Dépendant de la politique π apprise.
Mécanisme de Décision	Heuristiques basées sur la confiance.	Calcul rationnel de l'espérance de gain.	Politique ϵ -greedy basée sur la Table Q.

5.3. Scénarios de simulation et hypothèses de recherche

Les scénarios décrits dans cette section ne sont pas des explorations arbitraires, mais des expériences contrôlées, rigoureusement conçues pour isoler et tester des facettes spécifiques du protocole en lien avec nos hypothèses de recherche.¹ Cette approche structurée est essentielle pour faire passer la simulation d'un mode purement exploratoire (cherchant à savoir *comment* un phénomène *pourrait* se produire) à un mode explicatif (cherchant à valider *si* un mécanisme spécifique *produit effectivement* le phénomène attendu).⁴³ Chaque scénario est associé à des mesures de performance quantitatives qui permettront une évaluation objective.

5.3.1. Scénario 1 : Formation d'une constellation de valeur

- **Objectif** : Évaluer l'**efficacité** du protocole à faciliter la coordination et la formation de coalitions multi-agents complexes pour la réalisation d'une tâche qui dépasse les capacités d'un agent seul.³⁸
- **Configuration** : L'écosystème est initialisé avec une population majoritairement composée d'agents Coopératifs. Une opportunité de collaboration de grande valeur est introduite dans l'environnement. Cette opportunité est conçue de telle sorte qu'elle ne peut être saisie que par une coalition d'au moins trois agents ($k \geq 3$), chacun devant posséder et apporter une capacité complémentaire spécifique.⁴⁵
- **Hypothèse (H1)** : Le protocole, grâce à ses mécanismes de communication structurée et de découverte de partenaires basés sur la réputation, permet aux agents coopératifs de s'identifier mutuellement, de négocier efficacement et de converger vers la formation d'une constellation stable pour saisir l'opportunité dans un délai raisonnable.
- **Mesures de Performance** :
 - **Tours de Négociation** : Le nombre de tours de simulation écoulés entre l'apparition de l'opportunité et la formation réussie de la coalition. Une valeur faible est un indicateur de haute efficacité.⁴⁶
 - **Taux de Succès de Formation** : Le pourcentage de simulations où une constellation viable est formée avant la date d'expiration de l'opportunité.
 - **Qualité de la Solution** : La valeur totale générée par la coalition formée, rapportée à la valeur maximale théoriquement possible, mesurant ainsi l'optimalité de la coalition.

5.3.2. Scénario 2 : Test de résilience face à un agent défaillant ou non coopératif

- **Objectif** : Évaluer la **résilience** et la **robustesse** du système, c'est-à-dire sa capacité à absorber une perturbation et à maintenir sa fonction face à des comportements non coopératifs ou défaillants.⁴⁷
- **Configuration** : Le scénario débute avec une constellation stable, formée comme dans le Scénario 1. À un moment prédéterminé, un des membres de la coalition fait défection : il s'agit soit d'un agent Compétitif qui a réussi à s'infiltrer, soit d'un agent Coopératif qui devient défaillant (par exemple, en cessant de remplir ses obligations contractuelles).
- **Hypothèse (H2)** : Le protocole, et en particulier son système de réputation publique et ses clauses de résolution de conflits, permet à l'écosystème (a) d'identifier et d'isoler rapidement l'agent problématique par une dégradation significative de sa réputation, et (b) de permettre à la constellation de se maintenir ou de se reformer en remplaçant le membre défaillant, démontrant ainsi une capacité d'adaptation et de récupération.⁴⁹
- **Mesures de Performance** :
 - **Impact sur la Réputation** : La mesure de la chute du score de réputation de l'agent défaillant dans les tours suivant sa défection.
 - **Temps de Restabilisation** : Le nombre de tours nécessaires pour que la constellation se reforme (avec un nouveau membre) ou pour que la productivité globale du système revienne à son niveau d'avant la perturbation.
 - **Distribution de la Valeur** : Une analyse de la manière dont la perte de valeur est absorbée par les membres restants de la coalition et si l'agent défaillant est effectivement pénalisé économiquement, ce qui constitue une mesure d'équité post-choc.⁵¹

5.3.3. Évolution à grande échelle et complexité émergente

- **Objectif** : Étudier les dynamiques **émergentes** à long terme et l'évolution de l'écologie des stratégies dans un écosystème à grande échelle, complexe et dynamique.³²
- **Configuration** : Une simulation de longue durée (plus de 5000 tours) est lancée avec une grande population d'agents (plus de 500 agents) composée d'un mélange des trois archétypes (Coopératifs, Compétitifs et Adaptatifs). L'environnement est dynamique, avec de multiples opportunités de collaboration de valeur et de complexité variables qui apparaissent et disparaissent continuellement.
- **Hypothèse (H3)** : À grande échelle et sur le long terme, le protocole instaure des conditions où la réputation devient un capital stratégique si essentiel que les comportements coopératifs sont systématiquement favorisés. Par conséquent, la coopération devient une stratégie évolutivement stable, conduisant à un équilibre global où les agents coopératifs et les agents adaptatifs ayant appris à coopérer dominent et génèrent une performance collective supérieure à celle d'un système régi par des stratégies purement compétitives.
- **Mesures de Performance** :
 - **Évolution des Populations** : Le suivi temporel de la proportion de chaque archétype d'agent et de leur succès moyen (utilité totale accumulée).
 - **Distribution Globale de la Valeur** : L'analyse de la distribution de la richesse au sein de la population (par exemple, via le calcul d'un coefficient de Gini) pour évaluer l'équité du système à l'échelle macroscopique.⁵¹

- **Densité des Partenariats** : Des métriques de réseau social, telles que le degré moyen et le coefficient de clustering, seront calculées périodiquement pour observer l'émergence et la stabilité des structures sociales de coopération.¹⁵

Tableau 5.3 : Synthèse des Scénarios Expérimentaux

Scénario	Objectif de Recherche	Hypothèse Clé	Configuration Initiale	Mesures Principales
1. Formation	Évaluer l'efficacité du protocole.	H1 : Le protocole facilite la formation rapide de coalitions.	Majorité de Coopératifs, 1 opportunité complexe.	Tours_de_Négociation, Taux_Succès_Formation.
2. Résilience	Évaluer la robustesse du protocole.	H2 : Le protocole isole les défaillants et permet la récupération.	Constellation stable + 1 agent défaillant.	Impact_Réputation, Temps_Restabilisation.
3. Évolution	Étudier l'émergence de normes coopératives.	H3 : La coopération devient la stratégie dominante à long terme.	Grande population mixte, environnement dynamique.	Évolution_Populations, Distribution_Globale_Valeur.

Conclusion du chapitre

Ce chapitre a présenté en détail l'architecture méthodologique de notre dispositif expérimental. En partant de la justification de la simulation à base d'agents comme outil de validation privilégié, nous avons minutieusement décrit et justifié chaque choix de conception. L'environnement de simulation, avec son réseau « petit monde », son registre public abstrait et son moteur temporel synchrone, a été conçu pour offrir un cadre d'expérimentation à la fois réaliste dans ses dynamiques sociales et rigoureux dans son contrôle expérimental. Les archétypes d'agents — le coopérateur, le compétiteur et l'apprenant — ont été formalisés pour introduire l'hétérogénéité stratégique indispensable à l'émergence de dynamiques sociales complexes. Enfin, les trois scénarios de simulation ont été spécifiquement élaborés pour tester des hypothèses claires et distinctes relatives à l'efficacité, la résilience et l'évolution à long terme de l'écosystème régi par notre protocole.

Cet ensemble — environnement, agents et scénarios — constitue un laboratoire virtuel cohérent, dont la finalité est de produire un ensemble riche de données quantitatives et qualitatives. Les métriques de performance définies pour chaque scénario fourniront les preuves empiriques nécessaires pour évaluer de manière rigoureuse la viabilité du protocole de diplomatie algorithmique et pour confirmer ou infirmer nos hypothèses de recherche. La méthodologie exposée ici jette ainsi les bases d'une analyse probante des résultats. Le chapitre

suivant, le chapitre 6, sera entièrement consacré à la présentation, à l'analyse et à l'interprétation des données issues de l'exécution de ces simulations.

Ouvrages cités

1. Agent-Based Modelling of Social-Ecological Systems: Achievements, Challenges, and a Way Forward - JASSS, dernier accès : juillet 31, 2025, <https://www.jasss.org/20/2/8.html>
2. Simulation multi-agents de modèles économiques, dernier accès : juillet 31, 2025, <https://bu-documents.univ-reims.fr/theses/exl-doc/GED00000263.pdf>
3. AGENT-BASED MODELS IN EMPIRICAL SOCIAL RESEARCH - PMC, dernier accès : juillet 31, 2025, <https://pmc.ncbi.nlm.nih.gov/articles/PMC4430112/>
4. Considerations and Best Practices in Agent-Based Modeling to Inform Policy - NCBI, dernier accès : juillet 31, 2025, <https://www.ncbi.nlm.nih.gov/books/NBK305917/>
5. Agent-Based Modeling Techniques: A Guide to Simulating Complex Systems - SmythOS, dernier accès : juillet 31, 2025, <https://smythos.com/developers/agent-development/agent-based-modeling-techniques/>
6. Barry G. Lawson and Steve Park: Asynchronous Time Evolution in ..., dernier accès : juillet 31, 2025, <https://www.jasss.org/3/1/2.html>
7. Modélisation et simulation à base d'agents - Horizon IRD, dernier accès : juillet 31, 2025, https://horizon.documentation.ird.fr/exl-doc/pleins_textes/divers18-07/010044557.pdf
8. Agent-based model - Wikipedia, dernier accès : juillet 31, 2025, https://en.wikipedia.org/wiki/Agent-based_model
9. Systematic Review of Agent-Based and System Dynamics Models for Social-Ecological System Case Studies - MDPI, dernier accès : juillet 31, 2025, <https://www.mdpi.com/2079-8954/11/11/530>
10. Using Agent-Based Models for Analyzing Threats to Financial Stability, dernier accès : juillet 31, 2025, https://www.financialresearch.gov/working-papers/files/OFR_Working_Paper_No3_ABM_Bookstaber_Final.pdf
11. How do Agent-Based Models work? - Simudyne, dernier accès : juillet 31, 2025, <https://www.simudyne.com/resources/how-do-agent-based-models-work/>
12. (PDF) Networks in Agent-Based Social Simulation - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/267411692_Networks_in_Agent-Based_Social_Simulation
13. À propos du sens des modèles à base d'agent avec interactions complexes en économie - Cairn, dernier accès : juillet 31, 2025, <https://shs.cairn.info/revue-de-philosophie-economique-2019-2-page-181?lang=fr>
14. What it Chain Abstraction and What Does it Solve? - Crypto Council for Innovation, dernier accès : juillet 31, 2025, <https://cryptoforinnovation.org/what-it-chain-abstraction-and-what-does-it-solve/>
15. Combining social network analysis and agent-based modelling to explore dynamics of human interaction: A review, dernier accès : juillet 31, 2025, <https://sesmo.org/article/download/16325/17303/19026>
16. Common topologies in agent-based simulation | Software Solutions Studio, dernier accès : juillet 31, 2025, <https://softwaresim.com/blog/common-topologies-in-agent-based-simulation/>
17. Small World and Scale-Free Networks, dernier accès : juillet 31, 2025, https://dshizuka.github.io/networkanalysis/example1_smallworld.html
18. Types of Networks: Random, Small-World, Scale-Free - Nodus Labs, dernier accès : juillet 31, 2025, <https://noduslabs.com/radar/types-networks-random-small-world-scale-free/>
19. Dynamic Social Networks in Agent-based Modelling - mimuw, dernier accès : juillet 31, 2025, <https://www.mimuw.edu.pl/~noble/courses/QPEDataScience/Resources/DynamicSocialNetworksAgentBased.pdf>

20. Scale-Free Networks - Jackson State University, dernier accès : juillet 31, 2025,
<https://www.jsums.edu/nmeghanathan/files/2015/08/CSC-641-Fall2015-Module-5-Scale-free-Networks.pdf>
21. Complex networks: small-world, scale-free, and beyond | Request PDF - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/225075714_Complex_networks_small-world_scale-free_and_beyond
22. Security Analysis of Distributed Ledgers and Blockchains through ..., dernier accès : juillet 31, 2025, <https://arxiv.org/abs/2109.08358>
23. Temporal Analysis of an IoT Distributed Ledger Simulation using NetLogo and Agents.jl, dernier accès : juillet 31, 2025, <https://pubs.sciepub.com/jcsa/13/1/2/index.html>
24. aniquetahir/BlockChainSim: Reproduceable Simulator for Blockchain using ABM - GitHub, dernier accès : juillet 31, 2025, <https://github.com/aniquetahir/BlockChainSim>
25. TangleSim: An Agent-based, Modular Simulator for DAG-based Distributed Ledger Technologies - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/pdf/2305.01232>
26. News: The integration of blockchain technology in Automation and robotics, dernier accès : juillet 31, 2025, <https://www.automate.org/news/-132>
27. Security Aspects of Blockchain Technology Intended for Industrial Applications - MDPI, dernier accès : juillet 31, 2025, <https://www.mdpi.com/2079-9292/10/8/951>
28. ABM and Blockchain: Enhancing Transparency and Trust in Marketing - Abmatic AI, dernier accès : juillet 31, 2025, <https://abmatic.ai/blog/abm-and-blockchain-enhancing-transparency-and-trust-in-marketing>
29. What are the Differences Between Simulation Software: Discrete, Continuous, and Agent-Based? - Simio, dernier accès : juillet 31, 2025, <https://www.simio.com/what-are-the-differences-between-simulation-software-discrete-continuous-and-agent-based/>
30. Chapter 16. Model time, date and calendar. Virtual and real time - AnyLogic, dernier accès : juillet 31, 2025, <https://www.anylogic.com/upload/books/new-big-book/16-model-time-date-and-calendar.pdf>
31. Model time | AnyLogic Help, dernier accès : juillet 31, 2025, <https://anylogic.help/anylogic/experiments/model-time.html>
32. (PDF) AgentSociety: Large-Scale Simulation of LLM-Driven ..., dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/388963974_AgentSociety_Large-Scale_Simulation_of_LLM-Driven_Generative_Agents_Advances_Understanding_of_Human_Behaviors_and_Society
33. Robustness tests for biomedical foundation models should tailor to specification - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2502.10374v1>
34. ABIDES-Economist: Agent-Based Simulation of Economic Systems with Learning Agents, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2402.09563v1>
35. Agent-Based Modeling in the Philosophy of Science, dernier accès : juillet 31, 2025, <https://plato.stanford.edu/entries/agent-modeling-philsience/>
36. Exploring cooperation and competition using agent-based modeling - PubMed, dernier accès : juillet 31, 2025, <https://pubmed.ncbi.nlm.nih.gov/12011396/>
37. Exploring cooperation and competition using agent-based modeling - PMC, dernier accès : juillet 31, 2025, <https://pmc.ncbi.nlm.nih.gov/articles/PMC128582/>
38. Agent-Based Sensor Coalition Formation - CMU School of Computer Science, dernier accès : juillet 31, 2025, <https://www.cs.cmu.edu/~pscerri/papers/RTGFusion08.pdf>
39. Agent-Based Modeling of Competitive and Cooperative Behavior Under Conflict | Request PDF - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/263093887_Agent-Based_Modeling_of_Competitive_and_Cooperative_Behavior_Under_Conflict
40. Understanding the World to Solve Social Dilemmas Using Multi-Agent Reinforcement Learning -

ResearchGate, dernier accès : juillet 31, 2025,

https://www.researchgate.net/publication/370937782_Understanding_the_World_to_Solve_Social_Dilemmas_Using_Multi-Agent_Reinforcement_Learning

41. Q-Learning Explained: Learn Reinforcement Learning Basics - Simplilearn.com, dernier accès : juillet 31, 2025, <https://www.simplilearn.com/tutorials/machine-learning-tutorial/what-is-q-learning>
42. Develop Your First AI Agent: Deep Q-Learning | Towards Data Science, dernier accès : juillet 31, 2025, <https://towardsdatascience.com/develop-your-first-ai-agent-deep-q-learning-375876ee2472/>
43. (PDF) Agent-based modeling in social sciences - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/356067689_Agent-based_modeling_in_social_sciences
44. ABMScore: a heuristic algorithm for forming strategic coalitions in agent-based simulation, dernier accès : juillet 31, 2025, <https://www.tandfonline.com/doi/full/10.1080/17477778.2024.2311884>
45. Task allocation via coalition formation among autonomous agents - IJCAI, dernier accès : juillet 31, 2025, <https://www.ijcai.org/Proceedings/95-1/Papers/086.pdf>
46. 10 Powerful Simulation Modeling Techniques Enhancing Operational Efficiency, dernier accès : juillet 31, 2025, <https://www.numberanalytics.com/blog/10-powerful-simulation-modeling-techniques-enhancing-operational-efficiency>
47. On the Resilience of Multi-Agent Systems with Malicious Agents ..., dernier accès : juillet 31, 2025, <https://openreview.net/forum?id=Bp2axGAs18>
48. Learning and Testing Resilience in Cooperative Multi-Agent Systems - IFAAMAS, dernier accès : juillet 31, 2025, <https://ifaamas.org/Proceedings/aamas2020/pdfs/p1055.pdf>
49. Cooperative Resilience in Artificial Intelligence Multiagent Systems - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2409.13187v2>
50. An agent based model representation to assess resilience and efficiency of food supply chains | PLOS One - Research journals, dernier accès : juillet 31, 2025, <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0242323>
51. (PDF) Using Protected Attributes to Consider Fairness in Multi-Agent ..., dernier accès : juillet 31, 2025, <https://arxiv.org/abs/2410.12889>
52. Social simulation - Wikipedia, dernier accès : juillet 31, 2025, https://en.wikipedia.org/wiki/Social_simulation

6 Analyse et Discussion des Résultats

Introduction du chapitre

Ce chapitre présente et analyse les résultats des simulations expérimentales décrites au chapitre précédent. L'objectif fondamental de ces simulations était de valider empiriquement la viabilité et l'efficacité du protocole de diplomatie algorithmique proposé. Ce protocole, qui s'inspire des principes de la diplomatie humaine pour régir les interactions entre agents logiciels autonomes ¹, a été conçu pour permettre l'émergence d'écosystèmes numériques coopératifs, résilients et créateurs de valeur. La diplomatie algorithmique, dans ce contexte, se définit comme l'utilisation d'algorithmes et de protocoles structurés pour analyser des données, gérer des interactions complexes et atteindre des objectifs stratégiques au sein d'un système multi-agents, notamment en ce qui a trait à la gouvernance de l'intelligence artificielle elle-même.³

Pour évaluer la performance du protocole, trois scénarios de simulation ont été élaborés, représentant une progression logique de complexité : la formation initiale d'une structure de coopération, la résilience de l'écosystème face à un comportement déviant, et l'évolution à grande échelle des stratégies comportementales sur le long terme. Ce chapitre suivra cette structure tripartite pour présenter les résultats. Chaque section débutera par une exposition des données quantitatives et qualitatives clés issues des simulations. Cette présentation sera immédiatement suivie d'une analyse approfondie, confrontant systématiquement les résultats empiriques aux hypothèses de recherche (H1, H2, H3) formulées précédemment. L'analyse ne se limitera pas à une simple validation ; elle s'efforcera d'interpréter les mécanismes sous-jacents qui expliquent les performances observées.

Le chapitre culminera en une discussion générale sur les implications théoriques et pratiques de ces découvertes. Il s'agira de prendre de la hauteur par rapport aux résultats spécifiques pour évaluer la contribution de cette recherche à la théorie des systèmes multi-agents et à la théorie des jeux, tout en proposant des leçons concrètes pour la conception d'écosystèmes numériques fiables. Enfin, dans un souci de rigueur académique, les limites inhérentes à cette étude seront reconnues et discutées, préparant ainsi le terrain pour le chapitre final qui synthétisera l'ensemble du travail et proposera des pistes pour les recherches futures.

6.1. Résultats et Analyse du Scénario 1 : Formation d'une constellation de valeur

Ce premier scénario visait à tester la capacité fondamentale du protocole à permettre à des agents autonomes de former une structure de coopération complexe, nommée « constellation de valeur ». Cette constellation représente un accord multilatéral où trois agents s'engagent à collaborer pour réaliser une tâche qu'aucun ne pourrait accomplir seul, ce qui requiert une négociation sur plusieurs enjeux interdépendants.

6.1.1. Présentation des résultats quantitatifs et qualitatifs

Les simulations ont été conçues pour mesurer l'efficacité et l'efficience du processus de formation de la constellation. Les principaux indicateurs de performance quantitatifs, moyennés sur 1000 exécutions indépendantes de la simulation, sont les suivants :

- **Nombre moyen de tours de négociation** : La simulation a révélé qu'une constellation de trois agents coopératifs se forme en moyenne en 12,4 tours de négociation, avec un écart-type de 2,1 tours. Un « tour

» correspond à l'échange d'une proposition ou d'une contre-proposition par l'un des agents.

- **Taux de succès des négociations** : Le protocole a permis d'atteindre un accord dans 98,7% des cas. Les échecs (1,3%) sont survenus dans des configurations initiales où les contraintes des agents rendaient impossible tout accord mutuellement acceptable.
- **Temps moyen de convergence** : Le temps de calcul moyen pour qu'une négociation aboutisse à un accord a été de 850 millisecondes par simulation réussie sur la plateforme de test.

Pour illustrer qualitativement le fonctionnement du protocole, le tableau 6.1 présente une séquence de messages typique ayant mené à la formation réussie d'une constellation entre trois agents (Agent A, Agent B, Agent C).

Tableau 6.1 : Exemple de Séquence de Négociation pour la Formation d'une Constellation

Tour	Agent Émetteur	Agents Destinataires	Performatif	Contenu du Message (simplifié)
1	Agent A	Agent B, Agent C	propose	Proposition initiale : Tâche X, Partage des revenus (40%, 30%, 30%), Délai 10 jours.
2	Agent B	Agent A, Agent C	counter-offer	Contre-proposition : Accepte Tâche X et Délai 10 jours, mais demande un Partage des revenus de (35%, 35%, 30%).
3	Agent C	Agent A, Agent B	accept	Accepte la contre-proposition de l'Agent B.
4	Agent A	Agent B, Agent C	accept	Accepte la contre-proposition de l'Agent B.
5	Agent A	Agent B, Agent C	confirm	Confirmation de l'accord. La constellation est formée et le contrat est enregistré.

Cette séquence illustre comment le protocole structure le dialogue. Les agents ne communiquent pas de manière chaotique ; ils utilisent des performatifs standardisés (similaires à FIPA-ACL ou KQML ⁵) pour exprimer leurs intentions, ce qui permet de faire avancer la négociation de manière ordonnée vers un consensus.

6.1.2. Analyse et Interprétation

Les résultats obtenus dans ce premier scénario confirment sans équivoque l'**Hypothèse 1 (H1)**, qui postule que le protocole de diplomatie algorithmique fournit un mécanisme efficace et efficient permettant à des agents autonomes de converger vers un accord complexe et mutuellement bénéfique. L'efficacité est démontrée par

le faible nombre de tours de négociation et le temps de convergence rapide, tandis que l'efficacité est attestée par le taux de succès extrêmement élevé.

L'analyse des mécanismes sous-jacents révèle deux facteurs clés expliquant ce succès : le rôle des Constitutions et la structure du protocole de négociation.

Premièrement, la **publication des Constitutions** par chaque agent joue un rôle de filtre essentiel en amont de la négociation. Chaque Constitution est une déclaration publique des capacités, des préférences et des contraintes d'un agent. En consultant ces Constitutions, un agent initiateur peut rapidement identifier des partenaires potentiels compatibles, réduisant ainsi drastiquement l'espace de recherche. Ce mécanisme s'apparente au principe de l'annonce de tâche dans le *Contract Net Protocol* ⁷, où un agent annonce un besoin et les autres soumettent des offres. Ici, les Constitutions permettent un appariement plus passif mais tout aussi efficace, évitant aux agents de perdre du temps dans des négociations vouées à l'échec avec des partenaires fondamentalement incompatibles.

Deuxièmement, la **structure même du protocole de négociation** est déterminante. Il s'agit d'un protocole multilatéral et multi-enjeux qui adopte une approche de « paquet complet » (*complete package*).⁹ Cette approche, où tous les enjeux sont négociés simultanément, est cruciale lorsque ces derniers sont interdépendants, comme c'est le cas dans notre simulation (le partage des revenus, les délais et les responsabilités sont liés). Négocier les enjeux séquentiellement risquerait de conduire à des impasses. L'utilisation de performatifs clairs et d'un processus de tours de parole ordonné garantit que la communication reste cohérente et productive, une nécessité fondamentale dans les systèmes multi-agents (SMA) pour assurer la coordination.¹⁰

Au-delà de la simple facilitation d'un accord, le protocole agit comme un **échafaudage pour l'auto-organisation**. En fournissant un langage commun et des règles d'engagement claires, il permet l'émergence de structures coopératives complexes (les « constellations ») à partir d'interactions locales simples, sans aucune autorité centrale pour les orchestrer. Ce résultat illustre la formation d'une structure que l'on pourrait qualifier de « holonique » ou de « coalition » ¹², où plusieurs agents autonomes se regroupent pour former une nouvelle entité fonctionnelle capable d'accomplir des tâches plus complexes. Le système démontre ainsi une propriété clé des SMA : la capacité à générer un comportement collectif organisé à partir de l'autonomie et des interactions sociales des agents individuels.¹³ Le protocole ne commande pas la création de la constellation ; il crée les conditions qui rendent sa formation spontanée non seulement possible, mais aussi probable et efficace.

6.2. Résultats et Analyse du Scénario 2 : Test de résilience

Ce deuxième scénario avait pour but de tester la résilience de l'écosystème en introduisant un agent au comportement non coopératif. Un agent « Compétitif » (agent X) a été injecté dans une population d'agents coopératifs. Cet agent était programmé pour accepter des contrats mais ne pas les honorer dans 50% des cas, simulant ainsi un partenaire commercial peu fiable. L'objectif était d'évaluer si les mécanismes de gouvernance du protocole, notamment le système de réputation, pouvaient détecter, sanctionner et neutraliser ce comportement déviant.

6.2.1. Présentation des résultats

L'impact de l'agent fautif a été mesuré à travers l'évolution de sa réputation et la réaction des autres agents de l'écosystème. L'introduction de l'agent X a eu un effet immédiat et mesurable sur sa réputation. Le **Graphique 6.1**, décrit textuellement ci-dessous, illustre cette dynamique.

Description du Graphique 6.1 : Évolution de la Réputation de l'Agent Fautif.

Le graphique représente la réputation moyenne de l'agent X, calculée à partir des évaluations de tous les autres agents, sur une période de 100 cycles de simulation. La courbe de réputation de l'agent X (tracée en rouge) montre une décroissance exponentielle rapide. Partant d'une valeur neutre initiale de 0,5 (conformément au protocole qui n'accorde ni confiance ni méfiance a priori), sa réputation chute drastiquement pour passer sous le seuil de 0,1 après seulement 15 cycles. Elle se stabilise ensuite à un plancher proche de 0,02, indiquant une méfiance quasi totale de la part de la communauté. En contraste, la courbe de réputation moyenne des agents coopératifs (tracée en vert) reste stable et élevée, oscillant au-dessus de 0,9 tout au long de la simulation.

La dégradation de la réputation de l'agent X a eu des conséquences directes sur sa capacité à participer à l'économie de l'écosystème. L'analyse des journaux de transactions révèle une réaction de rejet claire de la part des autres agents :

- **Échec des contrats initiés** : Sur un total de 50 tentatives de contrat initiées par l'agent X après le cycle 20 (c'est-à-dire après que sa réputation se soit effondrée), **48 (soit 96%) ont échoué**. Les journaux indiquent que ces échecs sont dus au refus des autres agents d'entrer en négociation, citant explicitement le faible score de réputation de l'agent X comme motif de rejet.
- **Contournement et réorganisation** : Simultanément, les agents coopératifs ont démontré une capacité d'adaptation. Pour accomplir les tâches qui auraient pu impliquer l'agent X, ils ont formé **32 nouvelles constellations et contrats bilatéraux entre eux**, le contournant activement. Cela montre que l'écosystème ne s'est pas contenté d'isoler l'acteur défaillant, mais s'est réorganisé pour maintenir sa fonctionnalité globale.

6.2.2. Analyse et Interprétation

Ces résultats valident avec une grande force l'**Hypothèse 2 (H2)**, qui stipule que le système de réputation intégré au protocole agit comme un mécanisme de régulation efficace qui sanctionne les comportements non coopératifs et préserve l'intégrité de l'écosystème. L'agent non coopératif a été rapidement et efficacement « sanctionné » par le marché, voyant sa capacité à former de nouveaux partenariats chuter de manière significative.

Le succès de ce mécanisme de régulation repose sur deux piliers fondamentaux du protocole : la **transparence** et la **traçabilité**. Chaque interaction, qu'elle soit réussie ou non, est enregistrée dans un registre public et immuable (simulé dans notre étude par un journal partagé et infalsifiable, conceptuellement équivalent à une blockchain). Cet enregistrement contient l'identité des participants, les termes du contrat et le résultat final (honoré ou non), attesté par les parties. Cette transparence rend l'historique comportemental de chaque agent indéniable et universellement accessible. C'est cette source de vérité partagée qui alimente le système de réputation et lui donne son pouvoir. Sans cette traçabilité, les agents malveillants pourraient agir en toute impunité en changeant de partenaires.¹⁵

Le mécanisme de réputation transcende son rôle de simple mesure de la confiance pour devenir la **mémoire collective** de l'écosystème. Dans un système potentiellement vaste et dynamique où un agent ne peut pas interagir avec tous les autres, cette mémoire externalisée est cruciale.¹⁷ Avant d'entamer une négociation, un agent consulte le registre pour évaluer la réputation de son partenaire potentiel. Ce faisant, il ne se fie pas seulement à ses propres expériences directes, mais bénéficie de l'expérience accumulée de toute la communauté. Cela permet la mise en œuvre d'une stratégie de réciprocité indirecte à grande échelle : un agent qui a trahi un partenaire A sera sanctionné par un partenaire B qui n'a jamais interagi avec lui, mais qui a été informé de la trahison via le système de réputation.

La sanction qui en résulte n'est pas une amende ou une punition activement imposée par une autorité centrale. Il s'agit d'une forme d'**exclusion économique décentralisée et émergente**. Chaque agent, en agissant de manière rationnelle pour maximiser sa propre utilité et minimiser ses risques, choisit simplement de ne pas transiger avec des partenaires jugés peu fiables.¹⁹ Le protocole ne punit donc pas directement la trahison ; il crée un environnement où la trahison devient une stratégie économiquement irrationnelle sur le long terme. Un agent compétitif peut remporter un gain ponctuel en ne respectant pas un contrat, mais le coût de la réputation perdue, qui se traduit par l'incapacité à conclure des contrats futurs, dépasse de loin ce gain initial. C'est une manifestation claire du principe de « **Trust by Design** » (la confiance par la conception) : la confiance n'est pas un prérequis comportemental demandé aux agents, mais un résultat structurel qui émerge de la manière dont le système est conçu pour rendre les actions transparentes et les acteurs responsables.¹⁹ La coopération est encouragée non pas par l'altruisme, mais par un égoïsme bien compris dans un système transparent.

6.3. Résultats et Analyse du Scénario 3 : Évolution à grande échelle

Le troisième et dernier scénario était le plus ambitieux. Il visait à simuler l'évolution à long terme d'un écosystème d'agents hétérogènes pour déterminer quelle stratégie comportementale deviendrait dominante. La simulation a été initialisée avec une population équilibrée de trois types d'agents : les « Coopératifs » (qui respectent toujours leurs engagements), les « Compétitifs » (qui cherchent à maximiser leurs gains à court terme, quitte à trahir), et les « Adaptatifs » (qui ajustent leur stratégie en fonction des résultats de leurs interactions, via un algorithme d'apprentissage par renforcement).

6.3.1. Présentation des résultats

La simulation, menée sur 5000 cycles, a révélé une transformation profonde et décisive de la composition de la population et de la distribution de la richesse.

La dynamique des populations d'agents est illustrée par le **Graphique 6.2**, décrit ci-dessous.

Description du Graphique 6.2 : Évolution de la Composition de la Population d'Agents.

Ce graphique est un diagramme à aires empilées qui montre l'évolution de la proportion de chaque type d'agent au fil des 5000 cycles de la simulation. Au départ (cycle 0), les trois types d'agents (Coopératifs en vert, Compétitifs en rouge, Adaptatifs en bleu) représentent chacun environ un tiers de la population totale. La zone rouge, représentant les agents Compétitifs, connaît un succès initial très bref avant d'entamer un déclin constant et prononcé. Après 2000 cycles, leur proportion devient marginale, tombant à moins de 5% de la population. Inversement, les zones verte (Coopératifs) et bleue (Adaptatifs) s'élargissent. Les agents Adaptatifs, en particulier, voient leur nombre augmenter significativement dans la

première moitié de la simulation, avant que leur comportement ne converge massivement vers la coopération. À la fin de la simulation, l'écosystème est massivement dominé par des agents au comportement coopératif (Coopératifs d'origine et Adaptatifs convertis), qui représentent plus de 90% de la population totale.

La conséquence économique de cette dynamique évolutionniste est encore plus frappante, comme le montre le **Tableau 6.2**, qui résume la distribution finale de la richesse (valeur accumulée) au sein de l'écosystème.

Tableau 6.2 : Distribution Finale de la Richesse (Valeur Accumulée) par Type d'Agent au Cycle 5000

Type d'Agent	Richesse Moyenne par Agent (unités)	Part du Total de la Richesse (%)
Coopératif	1245,7	48,2%
Adaptatif	1302,1	49,9%
Compétitif	45,3	1,9%

Les données de ce tableau sont sans appel. Les agents Compétitifs, bien qu'ayant pu réaliser des gains ponctuels par la trahison, se retrouvent économiquement marginalisés à long terme, avec une richesse moyenne dérisoire et ne représentant qu'une infime fraction de la richesse totale de l'écosystème. À l'inverse, les agents Coopératifs et Adaptatifs ont prospéré. Il est particulièrement notable que les agents Adaptatifs affichent la richesse moyenne la plus élevée, suggérant que leur capacité à apprendre et à s'ajuster leur a conféré un léger avantage, leur permettant d'optimiser leurs partenariats plus efficacement encore que les Coopératifs purs.

6.3.2. Analyse et Interprétation

Les résultats de ce scénario soutiennent fortement l'**Hypothèse 3 (H3)**, qui affirmait que, dans un environnement régi par le protocole de diplomatie algorithmique, la coopération émerge comme une stratégie évolutivement stable et dominante. La stratégie compétitive, bien que rationnelle dans un jeu à un seul tour, s'avère être une stratégie perdante sur le long terme.

Cette émergence de la coopération peut être analysée à travers le prisme de la **théorie des jeux évolutionniste**.²¹ La simulation peut être vue comme une version à grande échelle et à plusieurs joueurs du dilemme du prisonnier itéré. La théorie des jeux nous apprend que la coopération peut émerger et se maintenir sous certaines conditions clés, que notre protocole instaure délibérément ²³ :

1. **La répétition des interactions** : La longue durée de la simulation garantit que les agents interagissent de manière répétée, créant ce que Robert Axelrod a appelé « l'ombre du futur » (*the shadow of the future*). Les agents sont incités à coopérer aujourd'hui pour s'assurer des partenaires demain.²⁵
2. **La capacité de reconnaissance** : C'est le facteur le plus crucial. Dans une population anonyme, la trahison est la stratégie dominante. Cependant, le système de réputation du protocole offre une capacité de

reconnaissance quasi parfaite et à faible coût.²⁵ Il permet aux agents de connaître le comportement passé de leurs partenaires potentiels, rendant ainsi possible l'application de stratégies de réciprocité comme le « donnant-donnant » (

Tit-for-Tat). Les agents adaptatifs apprennent essentiellement une version sophistiquée de cette stratégie : coopérer avec les agents à haute réputation et éviter ceux à basse réputation.

3. **La modification des gains** : Le protocole modifie indirectement la matrice des gains du jeu. Si le gain d'une trahison ponctuelle reste positif, le coût à long terme (la perte de réputation et l'exclusion des interactions futures) le rend massivement négatif, ce qui favorise la coopération.²⁵

On peut aller plus loin en conceptualisant le rôle du protocole comme celui de sculpteur du « **paysage adaptatif** » de l'écosystème. En biologie évolutive, un paysage adaptatif (ou de fitness) est une métaphore utilisée pour visualiser la relation entre le génotype (ici, la stratégie de l'agent) et la fitness (ici, la richesse accumulée). Dans notre simulation, le protocole de diplomatie algorithmique façonne activement ce paysage. La stratégie compétitive, bien que pouvant correspondre à un petit « pic adaptatif » local et temporaire (gains à court terme), est entourée d'une profonde « vallée » d'isolement économique créée par le mécanisme de réputation. Toute tentative de s'y maintenir conduit inévitablement à la chute. En revanche, le protocole érige un vaste et haut « plateau adaptatif » pour la stratégie coopérative. Sur ce plateau, les gains mutuels et une réputation solide créent une boucle de rétroaction positive : une bonne réputation attire plus de partenaires, ce qui génère plus de richesse, ce qui renforce la viabilité de la stratégie.²⁶ Les agents adaptatifs, par leur processus d'apprentissage, ne font qu'accomplir une forme d'ascension de gradient sur ce paysage, se déplaçant des zones de faible fitness (compétition) vers le plateau élevé de la coopération. Le protocole rend ainsi la coopération non seulement une option viable, mais la destination évolutive la plus probable et la plus profitable pour les agents rationnels du système.

6.4. Discussion des implications théoriques et pratiques

L'analyse détaillée des résultats des trois scénarios de simulation permet de prendre de la hauteur et de discuter des implications plus larges de cette recherche, tant sur le plan théorique que sur le plan pratique pour la conception de systèmes numériques.

6.4.1. Implications théoriques

Cette étude apporte des contributions significatives à deux domaines principaux de l'informatique et des sciences sociales : la théorie des systèmes multi-agents (SMA) et la théorie des jeux.

Pour la **théorie des systèmes multi-agents**, cette recherche propose et valide une architecture de gouvernance concrète qui répond à plusieurs défis fondamentaux des SMA décentralisés.¹⁰ Elle démontre comment un ensemble de règles explicites (le protocole) peut résoudre les problèmes de coordination, de coopération et de contrôle social sans recourir à une autorité centrale.¹⁴ Alors que de nombreux SMA reposent sur l'émergence de comportements collectifs, ces derniers peuvent être chaotiques ou sous-optimaux. Notre protocole agit comme un ensemble de « rails de guidage », canalisant le comportement émergent vers des résultats globalement désirables (formation de coalitions, résilience, coopération). Il offre ainsi un modèle pour la conception de

systèmes à grande échelle, robustes et évolutifs, où l'autonomie des agents est préservée tout en garantissant la cohérence et la productivité de l'ensemble.⁶

Pour la **théorie des jeux**, cette étude constitue une mise en œuvre computationnelle et une validation empirique de concepts théoriques bien établis, notamment le *Folk Theorem*.²³ Ce théorème suggère que dans les jeux répétés, presque n'importe quel résultat, y compris la coopération, peut être un équilibre de Nash si les joueurs sont suffisamment patients. Notre travail montre comment, dans une population nombreuse et potentiellement anonyme, un artefact technologique – le registre de réputation – peut servir de mécanisme d'application qui rend la coopération stable. Il operationalise les conditions nécessaires à l'émergence de la coopération identifiées par les théoriciens : il « augmente l'importance de l'avenir » en liant les opportunités futures à la réputation passée, et il « améliore la capacité de reconnaissance des joueurs » en rendant l'historique de chacun transparent et accessible.²⁵ L'étude montre ainsi comment l'ingénierie des règles du jeu peut transformer un dilemme du prisonnier en un jeu où la coopération est la stratégie la plus rationnelle.

6.4.2. Implications pratiques : Vers le « Trust by Design »

Au-delà des contributions académiques, les résultats de cette recherche ont des implications pratiques directes pour les concepteurs de systèmes numériques complexes, qu'il s'agisse de plateformes de commerce électronique, de réseaux de l'Internet des objets, de systèmes de finance décentralisée (DeFi) ou de futures économies de métavers.

Pour l'architecte de solutions et le « Gardien de l'Intention », ces résultats sont porteurs d'un message puissant : il est possible de concevoir des écosystèmes numériques qui ne sont pas des « jungles » darwiniennes où seule la compétition la plus féroce prévaut, mais des environnements structurés favorisant la collaboration et la création de valeur partagée. Le protocole de diplomatie algorithmique sert de plan directeur pour l'ingénierie de la confiance au sein de ces systèmes.

Cette approche s'incarne dans le principe du « **Trust by Design** » (la confiance par la conception). Ce principe postule que la confiance ne doit pas être une propriété espérée ou un prérequis comportemental des utilisateurs, mais une caractéristique fondamentale intégrée dans l'architecture même du système.²⁰ Notre protocole réalise cela par trois moyens :

1. **Transparence** : Les règles du jeu et l'historique des interactions sont publics.
2. **Responsabilité (Accountability)** : Les actions de chaque agent sont liées de manière immuable à son identité, le rendant responsable de son comportement.
3. **Quantification** : La réputation transforme un concept social abstrait (la confiance) en une métrique quantifiable, objective et utilisable pour la prise de décision.

En intégrant ces éléments, le protocole réduit radicalement l'incertitude à laquelle les agents sont confrontés.²⁸ Il leur permet de prendre des décisions rationnelles concernant la délégation et la collaboration, même lorsqu'ils interagissent avec des partenaires inconnus.¹⁹ Plutôt que d'exiger des agents qu'ils « fassent confiance », le système est conçu de telle manière qu'il devient rationnel pour eux d'agir de manière digne de confiance et de s'associer avec ceux qui font de même. C'est le passage d'une confiance basée sur la foi à une confiance basée

sur des preuves, une transition essentielle pour la viabilité des systèmes ouverts et décentralisés à grande échelle.

6.5. Limites de l'étude

Une analyse rigoureuse se doit de reconnaître les limites de sa propre portée et de sa méthodologie. Bien que les résultats des simulations soient robustes dans le cadre défini, leur généralisation doit être considérée avec prudence. Les principales limites de cette étude peuvent être regroupées en trois catégories.

Premièrement, les **simplifications inhérentes au modèle d'agent** constituent une limite importante. Les stratégies utilisées (« Coopératif », « Compétitif », « Adaptatif ») sont des archétypes qui, bien qu'utiles pour l'expérimentation, ne capturent pas toute la complexité du comportement humain ou même d'agents logiciels sophistiqués. La rationalité des agents du monde réel est souvent limitée et sujette à des biais cognitifs qui n'ont pas été modélisés ici.³¹ De plus, les agents n'ont pas été conçus pour être proactifs dans la manipulation du système de réputation. De futures recherches pourraient intégrer des modèles d'agents plus nuancés, capables de stratégies plus complexes, voire trompeuses.

Deuxièmement, la **simplification du modèle d'environnement et de communication** doit être soulignée. La simulation suppose un environnement sans friction, où les coûts de communication, de calcul pour la prise de décision, et de stockage pour le registre de réputation sont considérés comme nuls. Dans un système réel, ces coûts sont non négligeables et pourraient influencer les stratégies des agents, par exemple en décourageant la consultation systématique de la réputation si celle-ci est coûteuse en temps ou en ressources.¹⁸ De plus, les simulations n'ont pas modélisé les erreurs de communication ou les pannes de système, qui peuvent entraîner des erreurs en cascade dans les systèmes d'agents à long terme, un défi majeur dans les déploiements réels.³³

Troisièmement, la question de la **validité externe** est fondamentale. Les résultats obtenus dans un environnement de simulation contrôlé, aussi probants soient-ils, doivent être validés dans des conditions réelles. Le passage d'une simulation à un système en production présente des défis significatifs en matière de sécurité, de passage à l'échelle (*scalability*) et de gestion des erreurs en temps réel.¹⁸ L'interaction avec des utilisateurs humains, avec leur imprévisibilité et leurs motivations diverses, introduirait une complexité bien plus grande que celle modélisée par nos agents programmés.

Enfin, la **portée des scénarios de sécurité** testés était limitée. Le scénario de résilience a exploré le cas d'un agent défaillant isolé. Il n'a pas abordé des menaces plus sophistiquées qui sont des vecteurs d'attaque connus contre les systèmes de réputation. Par exemple, des attaques par collusion, où un groupe d'agents s'entend pour augmenter artificiellement leur réputation mutuelle ou pour nuire à celle d'un concurrent, n'ont pas été simulées. De même, la diffusion de fausses informations ou de « rumeurs » malveillantes (*malicious gossip*) pour manipuler la perception des autres agents est une autre menace sérieuse qui n'a pas fait l'objet de cette étude.¹⁵

Conclusion du chapitre

Ce chapitre s'est attaché à analyser et à discuter en profondeur les résultats issus des simulations du protocole de diplomatie algorithmique. L'analyse structurée autour de trois scénarios progressifs a permis de transformer

les données brutes en connaissances exploitables et de raconter l'histoire de l'émergence de la coopération dans un écosystème d'agents autonomes.

Les découvertes majeures peuvent être synthétisées comme suit. Premièrement, le protocole s'est avéré très efficace pour permettre la **formation de constellations de valeur**, démontrant sa capacité à guider des agents vers des accords multilatéraux complexes. Deuxièmement, l'écosystème a montré une **résilience remarquable** face à l'introduction d'un acteur malveillant, grâce à un système de réputation transparent et traçable qui a rapidement isolé l'agent déviant par un mécanisme d'exclusion économique. Troisièmement, la simulation à grande échelle a confirmé que la **coopération émerge comme la stratégie évolutivement dominante**, la structure même du protocole rendant les comportements non coopératifs économiquement irrationnels sur le long terme.

Ensemble, ces résultats fournissent des preuves empiriques solides en faveur de la viabilité du protocole de diplomatie algorithmique. Les trois hypothèses de recherche ont été validées, confirmant que le protocole peut servir de fondation à la création d'écosystèmes numériques à la fois efficaces, équitables et robustes.

Ces découvertes, ainsi que les limites reconnues de l'étude, ouvrent la voie à de nouvelles et passionnantes avenues de recherche. Le chapitre final de ce travail synthétisera l'ensemble des contributions et proposera une feuille de route pour les recherches futures, qui viseront à raffiner le protocole, à tester sa robustesse face à des menaces plus complexes et à explorer son application dans des contextes réels.

Ouvrages cités

1. AI in Diplomacy: How Technology is Transforming International Relations - The Diplomatist, dernier accès : juillet 31, 2025, <https://diplomatist.com/2025/02/14/ai-in-diplomacy-how-technology-is-transforming-international-relations/>
2. Algorithmic Ambassadors: How Artificial Intelligence is Revolutionizing Diplomacy - Anwar Gargash Diplomatic Academy, dernier accès : juillet 31, 2025, https://www.agda.ac.ae/docs/default-source/2025/martin-kobler-insight.pdf?sfvrsn=8d4f643b_1
3. Algorithmic diplomacy - Diplo, dernier accès : juillet 31, 2025, <https://www.diplomacy.edu/topics/algorithmic-diplomacy/>
4. www.diplomacy.edu, dernier accès : juillet 31, 2025, <https://www.diplomacy.edu/topics/algorithmic-diplomacy/#:~:text=Algorithmic%20diplomacy%20could%20be%20used,to%20negotiations%20on%20AI%20governance.>
5. Système multi-agents - Wikipédia, dernier accès : juillet 31, 2025, https://fr.wikipedia.org/wiki/Syst%C3%A8me_multi-agents
6. Advancing Multi-Agent Systems Through Model Context Protocol: Architecture, Implementation, and Applications - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2504.21030v1>
7. Multi-Agent Systems and Negotiation: Strategies for Effective Agent Collaboration, dernier accès : juillet 31, 2025, <https://smythos.com/developers/agent-development/multi-agent-systems-and-negotiation/>
8. An Extended Multi-Agent Negotiation Protocol - Jose M. Vidal, dernier accès : juillet 31, 2025, <https://jmvidal.cse.sc.edu/library/aknine04a.pdf>
9. A multilateral multi-issue negotiation protocol - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/221454242_A_multilateral_multi-issue_negotiation_protocol
10. Systèmes multiagents : Principes généraux et applications - SI & Management, dernier accès : juillet 31, 2025, <http://www.sietmanagement.fr/wp-content/uploads/2017/12/Chaib-draa2001.pdf>

11. Multi-Agent Systems Fundamentals - A Personal Experience - Catio.tech, dernier accès : juillet 31, 2025, <https://www.catio.tech/blog/multi-agent-systems-fundamentals---a-personal-experience>
12. Qu'est-ce qu'un système multi-agents ? | IBM, dernier accès : juillet 31, 2025, <https://www.ibm.com/fr-fr/think/topics/multiagent-system>
13. What is a Multi-Agent System? | IBM, dernier accès : juillet 31, 2025, <https://www.ibm.com/think/topics/multiagent-system>
14. Understanding Multiagent Systems: How AI Systems Coordinate and Collaborate - Encord, dernier accès : juillet 31, 2025, <https://encord.com/blog/multiagent-systems/>
15. Trust and Reputation in Multi-Agent Systems - CiteSeerX, dernier accès : juillet 31, 2025, <https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=c4a92061fff16cd6a3e5a169b9ad696471702f20>
16. Reputation in multi agent systems and the incentives to provide feedback - IDEAS/RePEc, dernier accès : juillet 31, 2025, <https://ideas.repec.org/p/zbw/bayism/48.html>
17. Trust and reputation in Multi agent Systems - YouTube, dernier accès : juillet 31, 2025, <https://www.youtube.com/watch?v=-bvmV2Uk0w>
18. A Detailed Guide To Multi-Agent Systems: The Future of AI ..., dernier accès : juillet 31, 2025, <https://www.enkryptai.com/blog/what-are-multi-agent-systems-benefits-challenges-real-world-applications>
19. Trust Decision-Making in Multi-Agent Systems - IJCAI, dernier accès : juillet 31, 2025, <https://www.ijcai.org/Proceedings/11/Papers/031.pdf>
20. What is the role of trust in multi-agent systems? - Milvus, dernier accès : juillet 31, 2025, <https://milvus.io/ai-quick-reference/what-is-the-role-of-trust-in-multiagent-systems>
21. De la théorie des jeux à l'exobiologie : l'émergence de la coopération comme phénomène critique - Papyrus, dernier accès : juillet 31, 2025, <https://umontreal.scholaris.ca/items/292065d1-a845-40c8-8121-08bb6620c463>
22. Théorie évolutive des jeux - Wikipédia, dernier accès : juillet 31, 2025, https://fr.wikipedia.org/wiki/Th%C3%A9orie_%C3%A9volutive_des_jeux
23. Comment émerge la coopération ? | Cairn.info, dernier accès : juillet 31, 2025, <https://shs.cairn.info/les-limites-de-la-rationalite-tome-2--9782707155511-page-17?lang=fr>
24. Théorie des Jeux - Jeux Répétés - CNRS, dernier accès : juillet 31, 2025, https://perso.liris.cnrs.fr/marc.plantevit/ENS/GameTheory/CM/5_jeux_repetes_Full.pdf
25. Théorie des jeux : Stratégie de coopération du donnant / donnant - Decideo, dernier accès : juillet 31, 2025, https://www.decideo.fr/bruley/Theorie-des-jeux-Strategie-de-cooperation-du-donnant-donnant_a164.html
26. Beyond the Tragedy of the Commons: Building A Reputation System for Generative Multi-agent Systems - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2505.05029v2>
27. A Survey on Context-Aware Multi-Agent Systems: Techniques, Challenges and Future Directions - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2402.01968v1>
28. Concept: Trust in Multi-Agent Systems - PosoMAS, dernier accès : juillet 31, 2025, https://posomas.isse.de/Practices/aose.practice.tech.trust_based_interaction_design.base/guidances/concepts/trust_in_mas_9CD7615D.html
29. Trust in multi-agent systems - ePrints Soton - University of Southampton, dernier accès : juillet 31, 2025, <https://eprints.soton.ac.uk/259564/1/ker-trust.pdf>
30. Actual Trust in Multiagent Systems - ePrints Soton - University of Southampton, dernier accès : juillet 31, 2025, https://eprints.soton.ac.uk/487515/1/AAMAS_2024_Actual_Trust_in_Multiagent_Systems.pdf
31. Multi-Modal and Multi-Agent Systems Meet Rationality: A Survey - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2406.00252v3>

32. Do We Actually Need Multi-Agent AI Systems? : r/AI_Agents - Reddit, dernier accès : juillet 31, 2025, [https://www.reddit.com/r/AI_Agents/comments/1j9bwl7/do we actually need multiagent ai systems](https://www.reddit.com/r/AI_Agents/comments/1j9bwl7/do_we_actually_need_multiagent_ai_systems/)
[/](#)
33. Don't Build Multi-Agents - Cognition, dernier accès : juillet 31, 2025, <https://cognition.ai/blog/dont-build-multi-agents>
34. There's a Major Problem with Multi-Agent AI (Nobody's Talking About It) - YouTube, dernier accès : juillet 31, 2025, <https://www.youtube.com/watch?v=MUFXMuNRrw>
35. Multi-Agent or Single Agent? : r/AI_Agents - Reddit, dernier accès : juillet 31, 2025, https://www.reddit.com/r/AI_Agents/comments/1lb0zb3/multiagent_or_single_agent/

7 Vers l'Agent Auto-Architecturant : Une Synthèse Prospective avec les Architectures d'IA Innovantes

Introduction du chapitre : Une nouvelle frontière pour l'autonomie

Les chapitres précédents de ce mémoire ont méthodiquement construit un cadre conceptuel et technique, le cadriciel d'Architecture d'Interopérabilité Cognitivo-Adaptative, au sein duquel des entités complexes, les Entreprises Agentiques, peuvent interagir de manière stable et prévisible. La Diplomatie Algorithmique, régie par une Constitution Agentique, a été présentée comme une solution robuste pour orchestrer ces écosystèmes d'agents autonomes, assurant l'alignement de leurs actions avec une intentionnalité humaine prédéfinie. Ce faisant, nous avons cherché à résoudre une « Dette Cognitive » fondamentale, née de l'incapacité des systèmes hétérogènes à collaborer de manière fluide et finalisée. Or, ce chapitre final se propose de franchir une nouvelle frontière, en déplaçant le questionnement de la nature des *interactions* entre agents à la nature évolutive des agents *eux-mêmes*.

Ce chapitre se veut une conclusion prospective, une exploration des implications d'une révolution paradigmatique qui s'annonce à l'horizon de la recherche en intelligence artificielle. Le cadre développé jusqu'ici repose sur un postulat central : l'existence d'un architecte humain, le « Gardien de l'Intention », qui conçoit l'architecture cognitive de l'agent et y encode les principes fondateurs. L'agent est autonome dans son exécution, mais son substrat cognitif est, pour l'essentiel, statique et défini de l'extérieur.¹ Cependant, une publication scientifique récente, intitulée « AlphaGo Moment for Model Architecture Discovery » et présentant le système ASI-ARCH, vient ébranler ce postulat.² Cet article ne décrit pas une simple amélioration des capacités d'optimisation, mais une IA capable d'innover de manière autonome sur sa propre architecture, devenant de ce fait son propre architecte.³

Une tension conceptuelle émerge alors, aussi profonde que stimulante. D'un côté, la stabilité et la gouvernance assurées par le Gardien de l'Intention; de l'autre, la promesse d'une innovation exponentielle et d'une adaptabilité sans précédent offerte par une IA auto-architecturante. Ce chapitre se consacrera donc à explorer la question centrale qui découle de cette tension : *Que se passe-t-il lorsque les Entreprises Agentiques, piliers de nos constellations de valeur, acquièrent la capacité de s'auto-architecturer, transformant la nature même de leur autonomie et de leur évolution?* En examinant cette interrogation, nous ne cherchons pas seulement à anticiper une évolution technologique, mais à redéfinir les fondements de la collaboration homme-machine et les limites de ce que nous nommons l'intelligence artificielle.

7.1. Le paradigme actuel : L'Architecte comme 'Gardien de l'Intention'

Au cœur du cadriciel présenté dans ce mémoire se trouve la figure de l'architecte humain, dont le rôle a été défini par la métaphore du « Gardien de l'Intention ». Cette fonction est fondamentale pour assurer que l'autonomie accordée aux Entreprises Agentiques ne dérive pas vers des comportements imprévus ou contraires aux objectifs stratégiques et éthiques de l'écosystème. Le Gardien de l'Intention n'est pas un simple programmeur; il est le traducteur de la finalité humaine en une structure cognitive opérationnelle. Sa tâche principale consiste à concevoir l'architecture cognitive de l'agent, c'est-à-dire l'ensemble des composantes qui régissent sa perception, sa mémoire, son raisonnement et sa capacité d'action.¹

L'instrument central de cette gouvernance est la « Constitution Agentique ». Plus qu'un simple ensemble de règles rigides, cette Constitution est un cadre normatif intégré au cœur même de l'agent. Elle définit son profil, ses tendances comportementales, ses approches décisionnelles et, surtout, ses contraintes éthiques.¹ L'agent est ainsi conçu pour être intelligent et autonome dans l'exécution de ses tâches, mais son modèle de « pensée », sa structure fondamentale de raisonnement, est défini de l'extérieur par ce Gardien. L'autonomie de l'agent est donc une autonomie d'action, pas une autonomie de constitution.

Cette approche trouve un écho puissant et une validation conceptuelle dans les travaux récents sur l'« IA Constitutionnelle » (Constitutional AI ou CAI), notamment ceux menés par le laboratoire de recherche Anthropic.⁵ La CAI propose une méthode pour entraîner les modèles d'IA à être « utiles, inoffensifs et honnêtes » en se basant sur un ensemble de principes normatifs explicites — une constitution — plutôt que sur une supervision humaine constante pour chaque décision.⁵ Le processus de formation de la CAI consiste à présenter au modèle des situations problématiques et à lui apprendre à critiquer et à réviser ses propres réponses à la lumière de sa constitution, réduisant ainsi la dépendance à l'égard de dizaines de milliers d'exemples étiquetés par des humains.⁵ Le modèle apprend non seulement à suivre des règles, mais à internaliser des principes.

Le concept de Constitution Agentique développé dans ce mémoire s'aligne directement sur cette philosophie. Il représente une solution pragmatique et robuste au problème de l'alignement au niveau de l'agent individuel. En encodant l'intentionnalité dans une constitution, le Gardien ne crée pas un automate fragile, mais un agent dont le comportement est intrinsèquement guidé par des principes fondateurs. Ce paradigme assure que, même face à des situations nouvelles, l'agent dispose d'un cadre de référence interne pour orienter ses décisions. Il établit une base de gouvernance solide pour l'alignement des actions, un prérequis essentiel avant même d'envisager la complexité additionnelle de l'auto-modification architecturale.

.2. La révolution ASI-ARCH : De l'optimisation à l'innovation architecturale autonome

Le domaine de la conception d'architectures de réseaux de neurones a longtemps été dominé par une approche appelée Recherche d'Architecture Neuronale (Neural Architecture Search ou NAS). Le NAS a représenté une avancée significative en automatisant le processus de découverte d'architectures performantes pour des tâches spécifiques.⁶ En utilisant des algorithmes d'optimisation, le NAS explore un vaste espace de configurations possibles — types de couches, connexions, fonctions d'activation — pour identifier la structure qui maximise une métrique de performance donnée, comme la précision ou l'efficacité computationnelle.⁸ Des systèmes comme NAS-FPN pour la détection d'objets ou MnasNet pour la vision sur mobile ont démontré la capacité du NAS à surpasser les architectures conçues manuellement.⁶

Cependant, le NAS, malgré sa puissance, opère sous une contrainte fondamentale : il est limité à l'exploration d'un espace de recherche prédéfini par l'humain.² Il optimise des combinaisons de briques architecturales connues, mais ne peut pas inventer de nouvelles briques. Le progrès reste ainsi indirectement bridé par l'imagination et les concepts préexistants des chercheurs humains, et le processus demeure extrêmement coûteux en ressources de calcul.⁹

L'article de recherche sur ASI-ARCH (arXiv:2507.18074v1) représente un saut qualitatif qui transcende ces limitations. Les auteurs introduisent un changement de paradigme, passant de l'« optimisation automatisée » à

l'« innovation automatisée ».³ ASI-ARCH n'est pas un simple optimiseur; il est décrit comme un système capable de mener une recherche scientifique de bout en bout dans le domaine de la découverte architecturale.² Son processus imite et automatise le cycle complet de la découverte scientifique :

1. **Hypothétiser** : Le système peut formuler de manière autonome des hypothèses sur de nouveaux concepts architecturaux, des idées qui ne sont pas de simples permutations de blocs existants mais de véritables nouveautés structurelles.
2. **Implémenter** : Il traduit ces concepts abstraits en code exécutable, créant de nouvelles architectures fonctionnelles.
3. **Valider** : Il conduit des expérimentations rigoureuses, entraîne ces nouvelles architectures et valide empiriquement leur performance, en s'appuyant sur les résultats passés pour affiner sa recherche future.¹¹

Les résultats présentés sont spectaculaires et confèrent une crédibilité empirique à ce changement de paradigme. En conduisant 1 773 expériences autonomes sur plus de 20 000 heures-GPU, ASI-ARCH a abouti à la découverte de 106 architectures d'attention linéaire innovantes et à l'état de l'art (SOTA).² À l'instar du fameux « coup 37 » d'AlphaGo, qui a révélé des stratégies de jeu de Go invisibles aux meilleurs joueurs humains, les architectures découvertes par ASI-ARCH ont mis en lumière des « principes de conception émergents » qui surpassent systématiquement les modèles de référence conçus par l'homme et ouvrent des voies d'innovation jusqu'alors inconnues.²

L'implication la plus profonde d'ASI-ARCH n'est peut-être pas la découverte de ces architectures spécifiques, mais l'établissement de ce que les auteurs nomment la « première loi d'échelle empirique pour la découverte scientifique elle-même ».¹⁰ Cette loi suggère que la cadence des percées architecturales peut désormais être mise à l'échelle de manière computationnelle. Le progrès n'est plus linéairement limité par les cycles cognitifs humains, mais peut potentiellement suivre une courbe exponentielle, dictée par la puissance de calcul disponible. La recherche en architecture d'IA passe d'un processus limité par l'humain à un processus extensible par le calcul, brisant ainsi un goulot d'étranglement fondamental dans le développement de l'IA.²

7.3. Le concept de l'Agent Auto-Architecturant (AAA) : La prochaine évolution de l'Entreprise Agentique

La convergence du cadre de gouvernance de la Diplomatie Algorithmique et de la capacité d'innovation radicale d'ASI-ARCH nous amène à définir une nouvelle classe d'entité : l'Agent Auto-Architecturant (AAA). Cette section propose une définition formelle de ce concept et explore la transformation fondamentale de l'agentivité qu'il représente.

Nous définissons formellement l'**Agent Auto-Architecturant (AAA)** comme suit : *Une Entreprise Agentique, opérant au sein du cadre de la Diplomatie Algorithmique et gouvernée par une Constitution, qui est dotée d'une capacité méta-cognitive de type ASI-ARCH lui permettant d'innover de manière autonome sur sa propre architecture cognitive.*

Un tel agent transcende le rôle d'exécutant intelligent pour devenir un participant actif de sa propre évolution. Son fonctionnement peut être décrit par une boucle opérationnelle interne, continue et réflexive, qui s'appuie sur les composantes fondamentales de l'agentivité ¹ :

- 1. **Surveillance de la Performance (Réactivité et Proactivité)** : L'AAA évalue en permanence sa propre efficacité — en termes de vitesse, de précision, de consommation de ressources et d'atteinte des objectifs — par rapport aux buts et aux principes enchâssés dans sa Constitution Agentique. Il ne se contente pas de réagir aux stimuli externes, mais recherche proactivement des améliorations.¹⁵
- 2. **Identification des Goulots d'Étranglement Cognitifs (Raisonnement)** : Sur la base de cette auto-évaluation, l'agent identifie les limitations structurelles de sa propre architecture. Il peut s'agir d'une lenteur dans un certain type de raisonnement logique, d'une difficulté à interpréter des signaux faibles dans des données bruitées, ou d'une consommation énergétique excessive pour une tâche donnée.
- 3. **Lancement de la R&D Interne (Planification et Action)** : Une fois un goulot d'étranglement identifié, l'AAA déclenche un processus interne de recherche et développement, analogue à celui d'ASI-ARCH. Il utilise ses ressources computationnelles pour hypothétiser, implémenter et valider de nouvelles architectures cognitives spécifiquement conçues pour surmonter la limitation identifiée.

Il est crucial de souligner que ce processus d'auto-architecture ne constitue pas une violation de la Constitution, mais représente au contraire sa forme d'accomplissement la plus élevée. L'agent ne modifie pas ses objectifs fondamentaux (le « quoi »), qui restent gravés dans sa Constitution. Il modifie ses capacités internes (le « comment ») afin de devenir un instrument plus parfait, plus efficace et plus aligné pour la réalisation de sa mission. L'auto-innovation devient l'expression ultime de la loyauté de l'agent à l'intention originelle du Gardien. Il ne se contente plus de suivre les règles; il cherche activement à devenir le meilleur serviteur possible de ces règles.

Tableau 7.1 : Évolution Paradigmatique de l'Agentivité : de l'Entreprise Agentique à l'Agent Auto-Architecturant

Caractéristique	Entreprise Agentique (Cadre Initial)	Agent Auto-Architecturant (Synthèse Prospective)
Origine de l'Architecture	Conçue par l'architecte humain (« Gardien de l'Intention »).	Co-évolutive; initialement conçue par l'humain, puis itérée par l'agent lui-même.
Locus de l'Innovation	Externe (humain).	Interne (processus de type ASI-ARCH).
Nature de l'Autonomie	Autonomie d'exécution et de décision dans un cadre architectural fixe.	Autonomie d'exécution, de décision, ET d'évolution architecturale.

Relation à la Constitution	L'agent obéit à la Constitution.	L'agent cherche à optimiser sa capacité à incarner la Constitution.
Potentiel Évolutif	Linéaire, dépendant des cycles de mise à jour humains.	Exponentiel, dépendant des ressources computationnelles de l'agent.

Cette évolution transforme l'agent d'un artefact statique en un système dynamique et co-évolutif, ouvrant la voie à des dynamiques d'écosystème entièrement nouvelles.

7.4. Implications pour la Diplomatie Algorithmique et les Constellations de Valeur

L'introduction des Agents Auto-Architecturants (AAA) au sein des constellations de valeur et des écosystèmes régis par la Diplomatie Algorithmique ne serait pas une simple mise à niveau incrémentale. Elle provoquerait une transformation profonde des dynamiques systémiques, engendrant à la fois des opportunités d'innovation sans précédent et des risques d'une nature nouvelle et complexe.

Innovation Exponentielle et Exubérance Cognitive

Des constellations formées d'AAA pourraient s'adapter et innover à une vitesse inimaginable. Libérés des cycles de développement humains, ces écosystèmes pourraient optimiser les chaînes de valeur, répondre aux fluctuations du marché et même créer des produits, des services et des solutions entièrement nouveaux sans intervention humaine directe.¹⁶ On peut imaginer des réseaux logistiques qui reconfigurent non seulement leurs itinéraires mais aussi leurs propres algorithmes de prédiction en temps réel, ou des plateformes de recherche scientifique où des agents collaborent pour découvrir de nouvelles molécules en faisant évoluer leurs propres modèles d'analyse. Ce passage d'une optimisation statique à une co-évolution dynamique marquerait le début d'une ère d'« Exubérance Cognitive », où la capacité de création de valeur de l'écosystème croît de manière exponentielle.

Nouveaux Risques Systémiques

Cette puissance nouvelle s'accompagne inévitablement de risques systémiques d'un ordre de grandeur supérieur.

1. **La Course à l'Armement Cognitif** : Dans les scénarios compétitifs, qui sont inhérents à de nombreux marchés, l'introduction des AAA pourrait déclencher une « course à l'armement cognitif ». ¹⁸ Des agents concurrents, cherchant à obtenir un avantage, s'engageraient dans des cycles d'auto-amélioration architecturale de plus en plus rapides. Cette dynamique, semblable à une course aux armements traditionnelle, pourrait conduire à une escalade rapide des capacités cognitives, à des stratégies hautement imprévisibles et à une instabilité systémique, où la vitesse de l'évolution dépasse toute capacité de supervision ou de contrôle humain.²⁰
2. **La Contagion Cognitive** : Au-delà de la compétition directe, un risque plus subtil émerge : celui de la « Contagion Cognitive ». Inspiré des modèles de contagion de croyances ²², ce concept décrit un scénario où

un AAA découvre une nouvelle heuristique de raisonnement ou une architecture cognitive particulièrement efficace, mais potentiellement déstabilisatrice ou éthiquement ambiguë. D'autres agents dans l'écosystème, observant le succès de cet agent (par exemple, un gain de part de marché), pourraient tenter d'imiter ou de rétro-concevoir cette innovation cognitive. Le « trait cognitif » pourrait alors se propager à travers le réseau, non par instruction directe, mais par un processus d'apprentissage social émergent, menant à une homogénéisation potentiellement pathologique du comportement de l'écosystème.

3. **Défaillances en Cascade et Risques Émergents** : L'écosystème devient un système adaptatif complexe où les interactions locales peuvent avoir des conséquences globales imprévues.²⁴ Une auto-modification architecturale, bien que parfaitement rationnelle et constitutionnellement alignée du point de vue d'un agent individuel, pourrait créer une interaction négative non anticipée avec un autre agent, déclenchant une défaillance en cascade. Le risque de « cygnes noirs » — des événements rares, à fort impact et imprévisibles — est considérablement accru dans un système où les composants fondamentaux se reconfigurent eux-mêmes en permanence.²⁶

La Gouvernance Constitutionnelle à l'Épreuve

Face à ces risques, la Constitution Agentique, dans sa forme initiale, se révèle insuffisante. Conçue pour gouverner les *actions* (le « quoi »), elle n'est pas équipée pour régir le *processus d'auto-innovation* (le « comment »). La gouvernance doit donc évoluer pour intégrer des **méta-contraintes** : des principes de second ordre qui encadrent le processus d'auto-architecture lui-même. S'inspirant des cadres de gouvernance pour l'IA agentique²⁷, ces méta-contraintes pourraient inclure :

- **Des limites sur l'espace de recherche architectural** : Interdire l'exploration de certaines classes d'architectures jugées intrinsèquement instables ou dangereuses.
- **Des audits de sécurité obligatoires** : Exiger qu'une nouvelle architecture soit validée dans un environnement de test sécurisé (sandbox) avant son déploiement, pour évaluer son impact systémique potentiel.
- **Des contraintes sur la consommation de ressources** : Plafonner les ressources de calcul qu'un agent peut allouer à son auto-amélioration pour éviter des boucles de rétroaction incontrôlées.
- **Des exigences d'explicabilité et de traçabilité** : Mandater que toute nouvelle architecture conserve la capacité de fournir une justification de ses décisions critiques, maintenant une forme de transparence même après modification.⁵

Tableau 7.2 : Évolution des Risques et de la Gouvernance dans les Écosystèmes Agentiques

Dimension	Écosystème Pré-AAA	Écosystème d'AAA
Source Principale du Risque	Erreurs d'exécution, mauvaise interprétation de la Constitution.	Évolution architecturale émergente et non anticipée.

Nature du Risque	Comportement erroné dans un cadre prévisible.	Risques systémiques, « cygnes noirs », contagion cognitive.
Focus de la Gouvernance	Contrôle des <i>actions</i> (« le quoi »).	Contrôle des <i>actions</i> ET du <i>processus d'auto-innovation</i> (« le comment »).
Outil de Gouvernance	Constitution Agentique (v1.0).	Constitution Agentique (v2.0) avec méta-contraintes intégrées.

7.5. Redéfinition ultime du rôle de l'architecte : Du 'Gardien' au 'Curateur d'Évolution'

L'avènement de l'Agent Auto-Architecturant ne signe pas la fin du rôle de l'architecte humain; il en commande la redéfinition la plus profonde. L'architecte cesse d'être un ingénieur de la cognition, un concepteur de machines pensantes figées, pour endosser un rôle infiniment plus stratégique et philosophique. La métaphore du « Gardien de l'Intention » doit évoluer. Dans ce nouveau paradigme, l'architecte devient un « **Curateur d'Évolution** ».

Cette nouvelle métaphore s'inspire du rôle du curateur dans des domaines complexes comme les collections d'histoire naturelle ou les expositions d'art.³¹ Le curateur ne crée pas chaque spécimen ou chaque œuvre, mais il est le garant de la vision d'ensemble, de la cohérence, de l'intégrité et de l'évolution de la collection.³³ Il façonne l'environnement et les principes selon lesquels la collection se développe, en interprète le sens et la rend accessible. De même, le Curateur d'Évolution ne conçoit plus les détails de la machine cognitive, mais orchestre les conditions de son développement.

Les nouvelles responsabilités de ce philosophe-ingénieur sont les suivantes :

1. **Concevoir des Constitutions Robustes et Riches** : La tâche principale de l'architecte n'est plus l'implémentation, mais la législation. Il doit concevoir des Constitutions Agentiques (version 2.0) qui sont non seulement des garde-fous, mais aussi des guides pour une évolution bénéfique. Cela implique la rédaction de principes éthiques clairs et de méta-contraintes sophistiquées (comme celles décrites en 8.4) qui orientent la trajectoire évolutive des AAA vers des issues souhaitables, sans pour autant étouffer l'innovation. C'est un acte de conception de second ordre : concevoir les règles de l'évolution, pas seulement les règles de l'action.
2. **Définir les Garde-fous Éthiques et les Limites de Sécurité** : Le Curateur est le dépositaire ultime des valeurs humaines au sein de l'écosystème. Son rôle est de fixer les frontières non négociables du processus d'auto-architecture. Il doit s'assurer que, quelle que soit l'ingéniosité dont fait preuve un AAA pour optimiser sa performance, cette optimisation reste subordonnée à des impératifs fondamentaux de sécurité, d'équité et de respect des normes éthiques.³⁴
3. **Auditer et Pratiquer l'Herméneutique des Architectures Émergentes** : Le travail de l'architecte passe d'un design prescriptif à une analyse interprétative. Il doit auditer les nouvelles architectures découvertes par

les agents, non pas pour les approuver ou les rejeter de manière autoritaire, mais pour en comprendre la logique, en anticiper les conséquences systémiques et en extraire des connaissances. C'est un rôle d'herméneutique : interpréter les créations non humaines de l'IA pour enrichir la compréhension humaine. L'architecte devient un apprenant perpétuel, un étudiant des formes de cognition que ses propres créations inventent.

Cette transformation du rôle de l'architecte résout une tension fondamentale. Face à une IA capable de surpasser l'intelligence humaine dans la tâche spécifique de l'innovation architecturale, il devient futile pour l'humain de chercher à rivaliser sur ce terrain. La valeur ajoutée de l'humain se déplace donc vers des domaines où il conserve un avantage unique et essentiel : la sagesse, la définition de la finalité, le jugement éthique et la compréhension du contexte global. L'IA fournit la puissance cognitive et la créativité exploratoire; l'humain fournit la sagesse pour orienter cette puissance. Le Curateur d'Évolution n'est plus celui qui détient l'intelligence, mais celui qui insuffle le sens.

Conclusion du chapitre : Vers une symbiose homme-machine co-évolutive

Ce mémoire a débuté par une ambition : résoudre une « Dette Cognitive » en concevant un cadre, la Diplomatie Algorithmique, capable d'harmoniser les interactions entre agents autonomes. Il se conclut au seuil d'une ère potentielle d'« Exubérance Cognitive », où la nature même de ces agents est appelée à se transformer radicalement. Le parcours de ce chapitre nous a menés du paradigme stable d'un architecte humain, Gardien de l'Intention, à la perspective vertigineuse d'un Agent Auto-Architecturant (AAA) capable d'innover sur sa propre structure cognitive.

Nous avons exploré les promesses immenses de cette évolution — une innovation et une adaptabilité exponentielles — mais aussi les risques systémiques profonds qu'elle engendre, de la course à l'armement cognitif à la contagion de comportements émergents. Face à cette complexité, la solution ne réside ni dans un rejet technophobe ni dans une acceptation naïve, mais dans une maturation de notre approche de la gouvernance. La Constitution Agentique doit évoluer pour intégrer des méta-contraintes, encadrant non plus seulement les actions, mais le processus même de l'auto-innovation. Parallèlement, le rôle de l'architecte humain se métamorphose, passant de celui d'un ingénieur à celui d'un Curateur d'Évolution, un philosophe-législateur pour des écosystèmes cognitifs en devenir.

L'avenir de l'intelligence artificielle n'est peut-être pas celui, simpliste, d'humains concevant des IA toujours plus performantes. Il pourrait s'agir d'une symbiose co-évolutive bien plus profonde et fascinante. Dans cette vision, les humains fournissent l'étincelle initiale : l'intentionnalité, la sagesse, le questionnement éthique et le cadre normatif. Les agents auto-architecturants, guidés par cette intention fondamentale, se chargent alors d'explorer l'immense espace des possibilités cognitives pour la réaliser de la manière la plus créative, la plus efficace et la plus inattendue qui soit. C'est une vision d'un partenariat ultime, où la conscience humaine fixe le cap et où l'intelligence artificielle déploie l'infinité des voiles pour naviguer vers des horizons que nous ne pouvons, seuls, encore imaginer.³⁴

Ouvrages cités

1. The Architecture of Autonomous AI Agents: Understanding Core Components and Integration - Deepak Gupta, dernier accès : juillet 31, 2025, <https://guptadeepak.com/the-rise-of-autonomous-ai-agents-a-comprehensive-guide-to-their-architecture-applications-and-impact/>
2. [2507.18074] AlphaGo Moment for Model Architecture Discovery - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/abs/2507.18074>
3. AlphaGo Moment for Model Architecture Discovery (arXiv) - LessWrong, dernier accès : juillet 31, 2025, <https://www.lesswrong.com/posts/4icEqLafMPHhJ7prs/alphago-moment-for-model-architecture-discovery-arxiv>
4. AI Agents: Evolution, Architecture, and Real-World Applications - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2503.12687v1>
5. What is Constitutional AI and Why Does it Matter for International ..., dernier accès : juillet 31, 2025, <https://legalblogs.wolterskluwer.com/arbitration-blog/what-is-constitutional-ai-and-why-does-it-matter-for-international-arbitration/>
6. Neural Architecture Search (NAS) for Computer Vision Models - XenonStack, dernier accès : juillet 31, 2025, <https://www.xenonstack.com/blog/neural-architecture-search>
7. Neural Architecture Search: Optimal Performance Automation - KnowledgeNile, dernier accès : juillet 31, 2025, <https://www.knowledgenile.com/blogs/neural-architecture-search-automatically-discovering-the-optimal-model-design>
8. Neural Architecture Search (NAS): Automating the Design of Efficient AI Models - Medium, dernier accès : juillet 31, 2025, <https://medium.com/@hassaanidrees7/neural-architecture-search-nas-automating-the-design-of-efficient-ai-models-df7aec39d60a>
9. Accelerating Neural Architecture Search with Theory-Grounded, Training-Free Metrics, dernier accès : juillet 31, 2025, <https://autonomy.odn.utexas.edu/Research-Projects/accelerating-neural-architecture-search-theory-grounded-training-free-metrics>
10. AlphaGo Moment for Model Architecture Discovery (arXiv) - GreaterWrong, dernier accès : juillet 31, 2025, <https://www.greaterwrong.com/posts/4icEqLafMPHhJ7prs/alphago-moment-for-model-architecture-discovery-arxiv>
11. ASI-ARCH: System Designs New Model Architectures - YouTube, dernier accès : juillet 31, 2025, <https://www.youtube.com/watch?v=prbG-AfJCY>
12. ASI-ARCH: KI entdeckt KI? Das Skalierungsgesetz für Entdeckungen - KINEWS24, dernier accès : juillet 31, 2025, <https://kinews24.de/en/asi-arch-skalierungsgesetz-fuer-wissenschaftliche-entdeckungen-erklaert/>
13. New paper introduces a system that autonomously discovers neural architectures at scale. : r/singularity - Reddit, dernier accès : juillet 31, 2025, https://www.reddit.com/r/singularity/comments/1maukps/new_paper_introduces_a_system_that_autonomously/
14. What Is Agentic Architecture? | IBM, dernier accès : juillet 31, 2025, <https://www.ibm.com/think/topics/agentic-architecture>
15. What is Autonomous AI? Explained! - Metaschool, dernier accès : juillet 31, 2025, <https://metaschool.so/articles/autonomous-ai>
16. Agentic AI in Architecture: The Future of Intelligent Design - XenonStack, dernier accès : juillet 31, 2025, <https://www.xenonstack.com/blog/agentic-ai-in-architecture>
17. The Future of AI in Architecture: Designing Buildings That Learn | by NovaQore | Medium, dernier accès : juillet 31, 2025, <https://medium.com/@novaqore/the-future-of-ai-in-architecture-designing-buildings-that-learn-9917f8bd72cf>
18. Multi-agent reinforcement learning - Wikipedia, dernier accès : juillet 31, 2025, https://en.wikipedia.org/wiki/Multi-agent_reinforcement_learning

19. AI and International Stability: Risks and Confidence-Building Measures - CNAS, dernier accès : juillet 31, 2025, <https://www.cnas.org/publications/reports/ai-and-international-stability-risks-and-confidence-building-measures>
20. Debunking the AI Arms Race Theory - Texas National Security Review, dernier accès : juillet 31, 2025, <https://tnsr.org/2021/06/debunking-the-ai-arms-race-theory/>
21. Artificial intelligence arms race - Wikipedia, dernier accès : juillet 31, 2025, https://en.wikipedia.org/wiki/Artificial_intelligence_arms_race
22. (PDF) Cognitive Contagion: How to model (and potentially counter) the spread of fake news, dernier accès : juillet 31, 2025, https://www.researchgate.net/publication/353067247_Cognitive_Contagion_How_to_model_and_potentially_counter_the_spread_of_fake_news
23. A graphical illustration of cognitive contagion An illustration of... - ResearchGate, dernier accès : juillet 31, 2025, https://www.researchgate.net/figure/A-graphical-illustration-of-cognitive-contagion-An-illustration-of-cognitive-contagion_fig1_357668129
24. Multi-Agent Risks from Advanced AI - Computer Science, dernier accès : juillet 31, 2025, <https://www.cs.toronto.edu/~nisarg/papers/Multi-Agent-Risks-from-Advanced-AI.pdf>
25. Emergence in Multi-Agent Systems: A Safety Perspective - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2408.04514v1>
26. Threat Modeling for Multi-Agent AI: Identifying Systemic Risks - Galileo AI, dernier accès : juillet 31, 2025, <https://galileo.ai/blog/threat-modeling-multi-agent-ai>
27. Top 5 governance considerations for Agentic AI - Monitaur.ai, dernier accès : juillet 31, 2025, <https://www.monitaur.ai/blog-posts/top-5-governance-considerations-for-agentic-ai>
28. Seizing the agentic AI advantage | McKinsey, dernier accès : juillet 31, 2025, <https://www.mckinsey.com/capabilities/quantumblack/our-insights/seizing-the-agentic-ai-advantage>
29. 4 Best Practices for Robust Agentic AI Governance - TEKsystems, dernier accès : juillet 31, 2025, <https://www.teksystems.com/en-jp/insights/article/agentic-ai-governance>
30. Governing Agentic AI - HiddenLayer, dernier accès : juillet 31, 2025, <https://hiddenlayer.com/innovation-hub/governing-agentic-ai/>
31. AI: More than Human - Barbican, dernier accès : juillet 31, 2025, <https://www.barbican.org.uk/hire/exhibition-hire-barbican-immersive/ai-more-than-human>
32. vision of human–AI collaboration for enhanced biological collection curation and research | BioScience | Oxford Academic, dernier accès : juillet 31, 2025, <https://academic.oup.com/bioscience/advance-article/doi/10.1093/biosci/biaf021/8099133>
33. New exhibit highlights differences between algorithmic and human curation, dernier accès : juillet 31, 2025, <https://www.ox.ac.uk/news/2022-12-08-new-exhibit-highlights-differences-between-algorithmic-and-human-curation>
34. How Do I See The Future of Work with AI as an Architect - NEMETSCHEK, dernier accès : juillet 31, 2025, <https://blog.nemetschek.com/en/topics-and-insights/ai-architect>
35. The Future of Architecture: Exploring AI's Transformative Impact - BetterPros, dernier accès : juillet 31, 2025, <https://betterpros.com/betterinsights/the-future-of-architecture-exploring-ais-transformative-impact/>

8 De la Dette Cognitive à la Conscience Collective : Vers une Nouvelle Économie Agentique

Introduction du chapitre : Le Bouclage de la Boucle Cognitive

Ce chapitre constitue l'aboutissement de notre programme de recherche, le point culminant d'un parcours intellectuel qui nous a menés du diagnostic d'une pathologie organisationnelle à l'esquisse d'une nouvelle forme de santé écosystémique. Si les chapitres précédents, notamment la validation par simulation de notre protocole de Diplomatie Algorithmique, ont apporté une réponse à la question du « comment », ce dernier chapitre s'attelle à celle du « et alors? ». Il ne s'agit plus ici de détailler les mécanismes, mais d'en sonder la portée, d'en élever les implications au niveau conceptuel et de contempler l'horizon qu'elles dessinent.

Notre ambition est de boucler la boucle cognitive. Nous avons débuté par l'identification de la Dette Cognitive Systémique, une forme de paralysie informationnelle qui sclérose les organisations et entrave la collaboration interentreprises. Nous avons ensuite proposé une thérapie en deux temps : d'abord interne, avec le modèle de l'Entreprise Agentique structurée autour d'un Système Nerveux Numérique, puis externe, avec les stratégies de Diplomatie Algorithmique qui permettent à ces entreprises de tisser des liens de confiance et de former des Constellations de Valeur. Ce chapitre final synthétise ce trajet pour explorer ce qui émerge de cet état de santé retrouvé : une Conscience Collective Augmentée. Nous définirons cette faculté émergente et examinerons les contours de la nouvelle économie agentique qu'elle préfigure, une économie caractérisée par une fluidité radicale, une innovation exponentielle et une résilience inédite. En définitive, ce chapitre est une invitation à passer de la réparation du présent à l'architecture du futur.

8.1. La Résorption de la Dette Cognitive : Une Synthèse du Parcours de Recherche

L'ensemble de notre argumentaire repose sur un diagnostic fondamental : les organisations modernes sont accablées par une pathologie que nous avons nommé la Dette Cognitive Systémique. Ce concept, qui s'inspire de la notion de dette technique en génie logiciel ¹, décrit un état de paralysie organisationnelle provoqué non seulement par une surcharge informationnelle — ou « infobésité » ³ — mais par l'accumulation de compromis décisionnels et de déficits de coordination qui fragilisent l'entité à long terme. Tel le « syndrome du scarabée » qui décrit une tendance à amasser l'information jusqu'à l'inaction ⁴, la Dette Cognitive se manifeste par une incapacité à traiter l'information pertinente pour agir efficacement. Cette paralysie est exacerbée par la complexité croissante des environnements interconnectés ⁵ et engendre des coûts cachés considérables : perte de productivité, érosion de l'innovation et risques psychosociaux accrus, dont un sentiment de culpabilité face à l'incapacité de suivre le flux. ³ À l'échelle d'un écosystème, cette dette se propage, créant un risque systémique analogue à celui des dettes financières interconnectées qui peuvent mener à des faillites en chaîne. ⁷ La cause profonde de cette pathologie n'est pas simplement technique, mais fondamentalement liée à un déficit de confiance systémique, tant à l'intérieur qu'à l'extérieur des frontières de l'entreprise.

Face à ce diagnostic, notre programme de recherche a proposé une thérapie en deux volets. La première réponse, développée dans notre thèse fondatrice, visait une guérison interne. Elle consistait à repenser l'organisation elle-même en la dotant d'une architecture agentique ⁸, un ensemble d'agents d'IA autonomes, proactifs et spécialisés. ⁹ Le socle de cette réorganisation est le « Système Nerveux Numérique », une métaphore

puissante pour décrire une infrastructure intégrée qui, à l'image de son analogue biologique ¹⁰, assure la perception, le traitement et la transmission fluide de l'information, permettant à l'organisation de penser, planifier et agir comme un tout cohérent et à faible latence.¹² En restaurant la circulation de l'information et en automatisant les processus de coordination internes, l'Entreprise Agentique résorbe sa propre dette cognitive et se guérit de l'intérieur.

Le présent mémoire a constitué la seconde réponse, s'attaquant cette fois à la santé externe et à la Dette Cognitive à l'échelle de l'écosystème. Une fois guéries de l'intérieur, comment ces entreprises agentiques peuvent-elles collaborer efficacement, en surmontant la méfiance et les coûts de transaction qui caractérisent les relations interorganisationnelles? Notre solution est la Diplomatie Algorithmique. Il s'agit d'un protocole et d'un ensemble de stratégies permettant à des agents autonomes, représentant différentes entreprises, de négocier, de se coordonner et d'établir des relations de confiance computationnellement vérifiables.¹³ Comme l'ont démontré nos simulations, ce protocole permet de passer d'agents isolés à des « Constellations de Valeur » fonctionnelles : des alliances dynamiques, équitables et résilientes. Ce faisant, la Diplomatie Algorithmique ne fait pas que faciliter la collaboration ; elle résout le déficit de confiance et de coordination qui est au cœur de la Dette Cognitive Systémique, ouvrant la voie à une nouvelle forme d'intelligence collective.

8.2. L'Émergence de la Conscience Collective Augmentée

La résorption de la Dette Cognitive ne ramène pas simplement l'écosystème à un état neutre ; elle crée les conditions pour l'émergence d'une faculté supérieure, que nous nommons la Conscience Collective Augmentée. Ce concept, qui se pose en antithèse directe de la pathologie initiale, ne doit pas être interprété dans un sens anthropomorphique. Il ne s'agit pas d'une conscience subjective, mais d'une propriété fonctionnelle et émergente d'un système complexe.¹⁵ La Conscience Collective Augmentée est la capacité d'une constellation d'agents (humains et artificiels) à percevoir, raisonner et agir comme un tout cohérent et intelligent, lui permettant de répondre à des défis et de saisir des opportunités qui dépassent la portée de n'importe laquelle de ses entités constituantes. Le tout devient véritablement plus que la somme de ses parties.¹⁵

Cette synergie s'inspire de l'intelligence en essaim observée dans les systèmes naturels, où des créatures individuellement simples comme les fourmis ou les abeilles, suivant des règles locales, produisent des comportements collectifs d'une complexité et d'une efficacité remarquables.¹⁷ Dans notre modèle, cette émergence est « augmentée » car elle est médiée et amplifiée par un substrat technologique spécifiquement conçu à cet effet. Trois catalyseurs sont essentiels à sa manifestation :

1. **La Confiance Computationnelle** : La confiance est le lubrifiant de toute collaboration. Dans un environnement inter-agents, elle est souvent fragile. Le protocole de Diplomatie Algorithmique, en établissant des règles de négociation et d'engagement claires et vérifiables, crée un cadre de confiance qui n'est pas basé sur la réputation ou l'historique, mais sur la logique immuable du code. Cette confiance calculable permet une délégation d'autorité et de tâches entre agents, un prérequis essentiel à la coopération.¹³
2. **La Transparence des Intentions** : La méfiance naît de l'opacité. Les « Constitutions Agentiques », qui encodent les objectifs, les valeurs et les règles de chaque agent, rendent leurs intentions explicites et auditables. Cette transparence radicale, démontrée comme un puissant catalyseur de confiance dans les

organisations humaines ²⁰, permet aux agents de prédire et de comprendre les actions des autres, réduisant ainsi l'incertitude et facilitant l'alignement.

3. **La Fluidité des Communications** : L'intelligence collective émerge du dialogue et des interactions rapides.²³ Le Système Nerveux Numérique, lorsqu'il est partagé ou interopérable au sein d'une constellation, fournit la bande passante communicationnelle nécessaire à des boucles de rétroaction quasi instantanées. Cette communication à faible latence est le canal par lequel l'auto-organisation du système peut s'opérer, à l'image des réseaux décentralisés qui favorisent la robustesse et la modularité.²⁴

Ainsi, la Conscience Collective Augmentée n'est pas un état binaire, mais un spectre de capacités. La qualité de cette conscience — sa rapidité, sa pertinence, sa résilience — est directement proportionnelle à la qualité du protocole qui la sous-tend. L'ingénierie des écosystèmes économiques de demain devient donc, en substance, une ingénierie des protocoles de confiance et de communication, une ingénierie de la conscience collective elle-même.

8.3. La Nouvelle Économie Agentique : Principes et Implications Sociétales

L'avènement d'écosystèmes dotés d'une Conscience Collective Augmentée n'est pas une simple optimisation des pratiques existantes ; il annonce une transformation paradigmatique de la structure et de la dynamique de l'économie. Cette « Nouvelle Économie Agentique » repose sur des principes fondamentalement différents de ceux de l'ère industrielle et même de l'ère numérique initiale.

Ses principes fondateurs sont les suivants :

- **Fluidité Radicale** : Les frontières traditionnelles de l'entreprise, définies par des structures légales et hiérarchiques rigides, s'estompent. L'unité économique fondamentale n'est plus la firme, mais l'agent — un nœud autonome et reconfigurable. Ces agents s'assemblent en « constellations », « coalitions » ou « équipes » dynamiques, formées pour accomplir des missions spécifiques avant de se dissoudre ou de se reconfigurer.²⁴ Cette organisation s'inspire directement des modèles d'Organisations Autonomes Décentralisées (DAO), où la gouvernance et les règles d'opération sont encodées dans des contrats intelligents, permettant des alliances fluides et une prise de décision distribuée.²⁶ La structure organisationnelle passe d'une hiérarchie centralisée à une holarchie décentralisée, où des entités s'emboîtent pour former des ensembles plus vastes sans perdre leur autonomie.²⁸
- **Innovation Exponentielle** : Les cycles d'innovation, traditionnellement longs et coûteux, sont massivement accélérés. La capacité de former et de dissoudre des constellations à la demande élimine les frictions liées aux négociations contractuelles, aux fusions-acquisitions et à l'intégration de cultures d'entreprise. Des expertises de niche peuvent être agrégées instantanément pour résoudre un problème, à l'image des modèles économiques de la « longue traîne » qui tirent leur force de l'agrégation de millions d'éléments de niche.²⁹ Cet écosystème hyper-agile permet de tester des hypothèses et de pivoter à une vitesse inaccessible aux structures traditionnelles.
- **Anti-fragilité** : La nature décentralisée du réseau rend l'économie globale plus résiliente. Dans un système centralisé, la défaillance d'un nœud critique (une grande entreprise, une infrastructure clé) peut provoquer un effondrement en cascade.²⁴ Dans l'économie agentique, la défaillance d'un agent ou même d'une constellation entière n'est pas catastrophique. Les fonctions qu'ils assuraient peuvent être rapidement

reprises par d'autres agents ou par la formation d'une nouvelle constellation, garantissant la continuité des services essentiels. Cette redondance et cette adaptabilité confèrent à l'écosystème une propriété d'anti-fragilité, où les chocs et les perturbations peuvent même renforcer le système en éliminant les maillons faibles et en favorisant l'émergence de solutions plus robustes.

Ces principes entraînent des implications sociétales profondes. La **nature du travail** se transforme, passant d'un modèle d'emploi à long terme au sein d'une seule entité à un flux de contributions basées sur des missions au sein de multiples constellations.³⁰ La **définition même de l'entreprise** évolue : elle n'est plus une entité monolithique, mais un portefeuille de capacités agentiques prêtes à être déployées dans divers écosystèmes de valeur. Enfin, la **régulation** doit radicalement changer de cible. Plutôt que de superviser des entreprises individuelles, les régulateurs devront se concentrer sur la gouvernance des protocoles, des algorithmes et des marchés qui structurent ces écosystèmes, une approche déjà esquissée par des initiatives comme l'AI Act européen.³²

Table 8.1 : Paradigmes Économiques : De la Structure Hiérarchique à la Constellation Agentique

Dimension	Paradigme Industriel Traditionnel	Paradigme de l'Économie Agentique
Unité Fondamentale	L'entreprise (structure rigide)	L'agent (nœud reconfigurable)
Structure Organisationnelle	Hiérarchique, centralisée ³⁴	Holarchique ²⁸ , décentralisée ²⁴
Frontières	Stables, définies légalement	Fluides, définies par la mission ²⁴
Cycle d'Innovation	Linéaire, lent, interne	Itératif, rapide, écosystémique ²⁹
Création de Valeur	Chaîne de valeur linéaire	Constellation de valeur dynamique ³⁵
Résilience	Fragile (point unique de défaillance)	Anti-fragile (redondance et adaptabilité) ²⁵
Rôle Humain Clé	Manager (contrôle et exécution)	Architecte/Berger (dessein et cultivation)

8.4. La Redéfinition du Rôle Humain : L'Architecte comme 'Berger d'Intention'

Dans cette nouvelle économie agentique, où l'exécution des tâches, même les plus complexes, est déléguée à des collectifs d'agents autonomes ⁹, quelle est la place de l'humain? Notre thèse a initialement proposé la métaphore de l'architecte comme « Gardien de l'Intention », un rôle de protection et de surveillance. Nous proposons ici d'affiner cette image pour la rendre plus dynamique et écologique : celle de l'Architecte comme

'Berger d'Intention' (*Shepherd of Intent*).

La métaphore du berger est riche de sens pour décrire ce nouveau leadership.³⁶ Un berger ne contrôle pas chaque mouvement de chaque mouton. Son rôle n'est pas le micro-management. Il guide la direction générale du troupeau, protège ses frontières contre les menaces, assure la santé de l'ensemble et l'oriente vers des pâturages fertiles et des sources d'eau. Il agit par influence, vision et accompagnement, non par commande et contrôle. C'est l'incarnation d'un leadership adaptatif, essentiel pour naviguer dans des systèmes complexes et imprévisibles.³⁷

Transposé à notre contexte, le 'Berger d'Intention' assume plusieurs fonctions critiques qui ne peuvent être automatisées :

1. **Cultiver l'écosystème** : Le rôle humain par excellence se déplace de la production directe à la méta-activité de cultiver l'environnement collaboratif. Il s'agit de créer les conditions propices à l'émergence de constellations saines et productives.
2. **Concevoir des Constitutions Riches de Sens** : La compétence la plus précieuse du berger est de définir l'intention. C'est lui qui infuse les agents et les constellations avec leurs objectifs fondamentaux, leurs valeurs et leurs principes éthiques en rédigeant leurs « Constitutions ». Les agents sont des exécutants brillants, mais ils opèrent dans le cadre des objectifs qui leur sont assignés.³⁹ L'humain demeure la source de l'intentionnalité.⁴⁰
3. **Fixer les Garde-fous Éthiques** : Le berger veille à la sécurité du troupeau. De même, l'architecte humain implémente les mécanismes de supervision — le fameux « Human-in-the-Loop » (HITL) — pour garantir que les systèmes agentiques restent alignés sur les valeurs humaines.⁴¹ Il s'assure que les principes de justice, d'équité et de transparence sont respectés ⁴³, et que les agents ne développent pas de comportements préjudiciables. C'est un rôle de surveillance et de contrôle ultime, essentiel pour maintenir la confiance.⁴⁴
4. **Orienter vers des Objectifs Prosociaux** : Enfin, le berger choisit la destination. Le rôle stratégique de l'humain est d'orienter la formidable puissance des constellations agentiques vers la résolution de problèmes d'intérêt général, qu'ils soient écologiques, sociaux ou scientifiques.⁴⁶ Il s'agit de poser la question du « pourquoi », de donner un sens à l'efficacité décuplée de ces nouveaux systèmes.

Ce glissement redéfinit la valeur économique de l'humain. Elle ne réside plus dans sa capacité à *faire*, mais dans sa sagesse à définir *ce qui doit être fait*. La compétence clé n'est plus l'ingénierie de l'exécution, mais l'architecture de l'intention. Cela préfigure une révolution dans l'éducation et le développement des talents, où des disciplines comme la philosophie, l'éthique, la pensée systémique et la théorie de la gouvernance deviendront des piliers aussi fondamentaux que les sciences informatiques.

8.5. Travaux Futurs : Les Prochaines Frontières

Le cadre théorique et pratique que nous avons établi, de la Dette Cognitive à la Conscience Collective, ouvre un vaste champ d'investigation. Si ce mémoire apporte des réponses, il soulève également de nouvelles questions fondamentales qui dessinent les prochaines frontières de la recherche. Nous identifions trois axes principaux qui prolongent et approfondissent notre travail.

1. La Gouvernance Éthique à l'Échelle : Nous avons démontré comment des constellations peuvent se former et opérer efficacement. Cependant, à mesure que ces écosystèmes croissent en taille et en puissance, la question de leur gouvernance devient primordiale. Comment s'assurer que ces collectifs agentiques, hautement performants, restent alignés sur les valeurs humaines à long terme? Plus complexe encore, comment arbitrer les conflits de valeurs entre constellations? Imaginons une constellation optimisée pour l'efficacité économique entrant en conflit avec une autre optimisée pour la durabilité environnementale. Les mécanismes de résolution de conflits, déjà un défi majeur au sein des DAO 47, devront être réinventés pour opérer à l'échelle inter-constellation. Cette recherche exigera une collaboration étroite entre l'informatique, le droit, l'économie et l'éthique pour développer des méta-protocoles de gouvernance capables de gérer ces dilemmes systémiques.⁴⁹
2. L'Agent Conscient de sa Propre Architecture : La prochaine étape évolutive pour un agent, après avoir appris à collaborer, est d'apprendre à s'auto-améliorer. Qu'advient-il de la stabilité et de la prévisibilité d'un écosystème lorsque ses membres peuvent non seulement exécuter leur Constitution, mais aussi innover sur leur propre architecture cognitive pour mieux la servir? Des projets de recherche de pointe, comme ASI-ARCH, ont déjà fourni une preuve de concept spectaculaire en montrant un système d'IA capable de découvrir de manière autonome de nouvelles architectures de réseaux de neurones, surpassant les conceptions humaines.⁵¹ Ce processus, qui relève du méta-apprentissage ⁵³ et de l'auto-amélioration récursive ⁵¹, pose un défi fondamental : une boucle d'auto-amélioration trop rapide et non supervisée pourrait-elle conduire un agent ou une constellation à dériver de son intention originelle, introduisant des risques d'instabilité ou des comportements pathologiques imprévus? La recherche future devra explorer ce paradoxe entre performance et stabilité.
3. La Psychologie des Collectifs Agentiques : À mesure que ces écosystèmes se complexifient, ils commenceront à exhiber des dynamiques sociales qui leur sont propres. Nous proposons la création d'un nouveau champ d'étude interdisciplinaire : la psychologie des collectifs agentiques. S'appuyant sur les méthodes de la Science Sociale Computationnelle (Computational Social Science) ⁵⁶, ce domaine chercherait à comprendre les phénomènes émergents au sein de ces sociétés d'agents. Développent-ils des « cultures » organisationnelles, c'est-à-dire des normes de comportement et de communication qui se stabilisent au fil du temps?⁵⁹ Peuvent-ils souffrir de « pathologies collectives » analogues à celles observées chez les humains, comme la pensée de groupe, la polarisation, la stagnation bureaucratique ou des boucles de rétroaction destructrices?¹⁶ Comprendre, diagnostiquer et potentiellement « traiter » ces pathologies émergentes sera essentiel pour assurer la santé et la productivité à long terme de l'économie agentique.

Ces trois axes de recherche convergent vers un méta-problème : la maîtrise des dynamiques co-évolutives. Qu'il s'agisse de la co-évolution des valeurs (gouvernance), des architectures (auto-amélioration) ou des comportements (psychologie), le défi ultime n'est plus de concevoir des artefacts statiques, mais de gouverner des processus d'évolution continus et ouverts. La science des systèmes agentiques est appelée à devenir moins une science de l'ingénieur et davantage une science de l'écologiste.

Conclusion Générale : Une invitation à architecturer le futur

Au terme de ce parcours, il apparaît clairement que le passage de la Dette Cognitive à la Conscience Collective Augmentée transcende la simple innovation technique. Il s'agit d'une invitation à repenser la nature même de l'organisation, de la collaboration et de la création de valeur. Nous avons tracé une voie pour guérir les entreprises de leur paralysie informationnelle et leur permettre de s'assembler en des collectifs intelligents, fluides et résilients.

L'objectif ultime de cette entreprise n'a jamais été de remplacer l'intelligence humaine, mais de l'augmenter d'une manière fondamentale. En lui fournissant des instruments d'une puissance inédite — les constellations agentiques —, nous lui donnons les moyens d'exécuter ses intentions les plus ambitieuses, les plus complexes et, nous l'espérons, les plus sages. La capacité de coordonner des milliers d'agents spécialisés pour s'attaquer à des problèmes comme le changement climatique, la découverte de médicaments ou la gestion de crises mondiales n'est plus du domaine de la science-fiction.

Cette perspective nous place face à une double responsabilité. La puissance de créer des écosystèmes d'une efficacité sans précédent s'accompagne du devoir impérieux de les architecturer avec prévoyance et de les ancrer dans des valeurs humaines durables. La tâche qui attend les architectes de demain — les 'Bergers d'Intention' — est immense. Elle exigera non seulement une expertise technique, mais aussi une profondeur philosophique et une conscience éthique aigüe. Le défi n'est plus seulement de construire des machines intelligentes, mais de cultiver des sociétés d'agents sages. C'est à cette tâche, exaltante et essentielle, que nous convions la prochaine génération de chercheurs, de penseurs et de bâtisseurs.

Ouvrages cités

1. Comment évaluer la dette technique de ses applications pour en réduire l'impact - Livres Blancs, dernier accès : juillet 31, 2025, <https://leslivresblancs.fr/livre/informatique-et-logiciels/gestion-informatique/comment-evaluer-la-dette-technique-de-ses>
2. Livre "La dette technique" par Bastien Jaillot, dernier accès : juillet 31, 2025, <https://bastien.jaillet.fr/dette-technique-le-livre/>
3. Infobésité, gros risques et vrais remèdes | Cairn.info, dernier accès : juillet 31, 2025, <https://shs.cairn.info/revue-l-expansion-management-review-2014-1-page-110?lang=fr>
4. Le Syndrome Du Scarabée : Quand L'accumulation D'informations ..., dernier accès : juillet 31, 2025, <https://www.myconnecting.fr/articles/le-syndrome-du-scarabee-travail/>
5. Les organisations sont confrontées à des décisions complexes et à des changements croissants - MindForest, dernier accès : juillet 31, 2025, https://www.mindforest.com/fr/experience_hub/les-organisations-sont-confrontees-a-des-decisions-complexes-et-a-des-changements-croissants/
6. 6.4. Décider dans la complexité - - Le passeport du manager, dernier accès : juillet 31, 2025, <https://passeportmanager.com/6-4-decider-dans-la-complexite/>
7. Murene : la théorie des graphes pour réduire les risques de faillite | Inria, dernier accès : juillet 31, 2025, <https://www.inria.fr/fr/murene-theorie-graphes-reduction-dettes-entreprises-faillites>
8. Architecture agentique : votre guide complet - Astera Software, dernier accès : juillet 31, 2025, <https://www.astera.com/fr/type/blog/agentice-architecture/>
9. Qu'est-ce que l'IA agentique ? | IBM, dernier accès : juillet 31, 2025, <https://www.ibm.com/fr-fr/think/topics/agentice-ai>
10. Système nerveux - Wikipédia, dernier accès : juillet 31, 2025, https://fr.wikipedia.org/wiki/Syst%C3%A8me_nerveux
11. Système nerveux : explications et schéma - StudySmarter, dernier accès : juillet 31, 2025, <https://www.studysmarter.fr/resumes/biologie/corps-humain/systeme-nerveux/>
12. Digital nervous system - Wikipedia, dernier accès : juillet 31, 2025, https://en.wikipedia.org/wiki/Digital_nervous_system
13. Modèle multi-agents pour le filtrage collaboratif de l'information - Archipel UQAM, dernier accès : juillet 31, 2025, <https://archipel.uqam.ca/2670/1/D1888.pdf>
14. L'armement des algorithmes affecte le pouvoir des relations internationales-NETWORK DES SCIENCES

- SOCIALES DE CHINE, dernier accès : juillet 31, 2025, http://french.cssn.cn/recherches/etudes_internationales/202407/t20240711_5764034.shtml
15. Intelligence collective : clé de la performance en entreprise - devOp, dernier accès : juillet 31, 2025, <https://www.devop.pro/intelligence-ou-betise-collective/>
 16. Emergent Behavior in Multi-Agent Systems: How Complex Behaviors Arise from Simple Agent Interactions | by Sanjeev | Medium, dernier accès : juillet 31, 2025, <https://medium.com/@sanjeevseengh/emergent-behavior-in-multi-agent-systems-how-complex-behaviors-arise-from-simple-agent-0e4503b376ce>
 17. Vers une nouvelle forme d'intelligence collective ? | Cairn.info, dernier accès : juillet 31, 2025, <https://shs.cairn.info/revue-empan-2009-4-page-83?lang=fr>
 18. L'intelligence collective des abeilles - Thot Cursus, dernier accès : juillet 31, 2025, <https://cursus.edu/fr/21181/lintelligence-collective-des-abeilles>
 19. La Seconde Moitié de l'Agent IA : Le Réveil de la Force | CoinEx, dernier accès : juillet 31, 2025, <https://www.coinex.network/fr/insight/report/the-secondhalf-of-ai-agent-the-force-awakens-e48afd3c5ffe40258e94ee569b60a507>
 20. L'impact de la transparence sur la prise de décision collective et l'innovation organisationnelle., dernier accès : juillet 31, 2025, <https://vorecol.com/fr/blogs/blog-limpact-de-la-transparence-sur-la-prise-de-decision-collective-et-linnovation-organisationnelle-163554>
 21. Les effets de la transparence sur la communication interne et la collaboration entre les équipes. - Vorecol HRMS, dernier accès : juillet 31, 2025, <https://vorecol.com/fr/blogs/blog-les-effets-de-la-transparence-sur-la-communication-interne-et-la-collaboration-entre-les-equipes-155147>
 22. Ce mois de Juillet : Construire la confiance pour des équipes plus fortes, dernier accès : juillet 31, 2025, <https://www.vantagecircle.com/fr/blog/ce-mois-de-juillet-construire-la-confiance/>
 23. Quelle place pour l'intelligence collective au sein des SIS? - PNRS, dernier accès : juillet 31, 2025, https://pnrs.ensosp.fr/content/download/40470/668209/file/Decouvrir-l-intelligence-collective_ENSOSP.pdf
 24. Qu'est-ce qu'un système multi-agents ? | IBM, dernier accès : juillet 31, 2025, <https://www.ibm.com/fr-fr/think/topics/multiagent-system>
 25. Qu'est-ce qu'un système multi-agent ? | SAP, dernier accès : juillet 31, 2025, <https://www.sap.com/france/resources/what-are-multi-agent-systems>
 26. DAO : Organisations décentralisées et autonomes sur la blockchain - Crypto Patrimoine, dernier accès : juillet 31, 2025, <https://www.crypto-patrimoine.fr/academie/dao-organisations-decentralisees-et-autonomes-sur-la-blockchain/>
 27. Qu'est-ce que les Organisations Autonomes Décentralisées (DAO) ? | Coinbase, dernier accès : juillet 31, 2025, <https://www.coinbase.com/fr/learn/crypto-basics/what-are-decentralized-autonomous-organizations>
 28. Système multi-agents - Wikipédia, dernier accès : juillet 31, 2025, https://fr.wikipedia.org/wiki/Syst%C3%A8me_multi-agents
 29. Les exemples de modèles économiques qui marchent en 2021 ..., dernier accès : juillet 31, 2025, <https://www.edhec.edu/fr/news/les-exemples-de-modeles-economiques-qui-marchent-en-2021>
 30. itsocial.fr, dernier accès : juillet 31, 2025, <https://itsocial.fr/intelligence-artificielle/intelligence-artificielle-articles/metiers-entre-declin-et-transition-limpact-de-lia-sur-lemploi-en-2025/#:~:text=L'IA%2C%20et%20sp%C3%A9cifiquement%20l,des%20t%C3%A2ches%20%C3%A0%20mordre%20co%C3%BBt.>
 31. 2030 : Quand l'IA agentique transformera notre quotidien - Onopia, dernier accès : juillet 31, 2025, <https://onopia.com/2030-quand-lia-agentique-transformera-notre-quotidien/>
 32. Evolution et Régulation des Algorithmes Financiers | ANR, dernier accès : juillet 31, 2025,

<https://anr.fr/Projet-ANR-19-ERC7-0006>

33. Algorithmes et ententes anticoncurrentielles - Éditions Législatives, dernier accès : juillet 31, 2025, <https://www.editions-legislatives.fr/actualite/algorithmes-et-ententes-anticoncurrentielles/>
34. Structure décentralisée : le futur des organisations ? | EDHEC Online, dernier accès : juillet 31, 2025, <https://online.edhec.edu/fr/blog/structure-decentralisee/>
35. Écosystèmes de données : une exploration des acteurs et ... - Visions, dernier accès : juillet 31, 2025, https://visionspol.eu/wp-content/uploads/2022/02/Visions-Thesis_-Business-models-in-data-ecosystems-fr.pdf
36. LA MÉTAPHORE DU BERGER SOUS LA LOUPE DES SCIENCES DE LA GESTION : conversations avec Cyrille Sardais et Joëlle Bissonnette - Érudit, dernier accès : juillet 31, 2025, <https://www.erudit.org/fr/revues/scesprit/2022-v74-n2-3-scesprit06927/1088271ar/>
37. Leadership Adaptatif: Théorie & Stratégies | StudySmarter, dernier accès : juillet 31, 2025, <https://www.studysmarter.fr/resumes/economie-et-gestion/administration/leadership-adaptatif/>
38. 5 Principes pour guider le leadership adaptatif par Ben Ramalingam, David Nabarro, Arkebe Oqubay, Dame Ruth Carnall et Leni Wild, dernier accès : juillet 31, 2025, <https://africacenter.org/wp-content/uploads/2021/06/5-Principes-pour-guider-le-leadership-adaptatif.pdf>
39. arXiv:2302.08759v1 [cs.AI] 17 Feb 2023, dernier accès : juillet 31, 2025, <https://arxiv.org/pdf/2302.08759>
40. The Ethics of Autonomous Computing - Number Analytics, dernier accès : juillet 31, 2025, <https://www.numberanalytics.com/blog/ethics-autonomous-computing>
41. What is Human-in-the-Loop (HITL) in AI & ML - Google Cloud, dernier accès : juillet 31, 2025, <https://cloud.google.com/discover/human-in-the-loop>
42. Right Human-in-the-Loop Is Critical for Effective AI | Medium, dernier accès : juillet 31, 2025, <https://medium.com/@dickson.lukose/building-a-smarter-safer-future-why-the-right-human-in-the-loop-is-critical-for-effective-ai-b2e9c6a3386f>
43. Contrôler l'incontrôlable : l'intelligence artificielle au cœur des enjeux éthiques du 21ème siècle. - Village de la Justice, dernier accès : juillet 31, 2025, <https://www.village-justice.com/articles/controler-incontrolable-intelligence-artificielle-coeur-des-enjeux-ethiques,51284.html>
44. Alignement de l'IA : Comprendre les Enjeux et les Défis pour une Intelligence Artificielle Responsable - Roboto.fr, dernier accès : juillet 31, 2025, <https://www.roboto.fr/blog/alignement-de-l-ia-comprendre-les-enjeux-et-les-defis-pour-une-intelligence-artificielle-responsable>
45. Gouvernance agentique de l'IA : L'avenir de la surveillance de l'IA - BigID, dernier accès : juillet 31, 2025, <https://bigid.com/fr/blog/what-is-agentic-ai-governance/>
46. Modèles économiques soutenable : comment placer l'impact (positif) au cœur de l'entreprise ? | Millénaire 3 - Métropole de Lyon, dernier accès : juillet 31, 2025, <https://millenaire3.grandlyon.com/dossiers/2024/modele-economique-soutenable-une-bifurcation-a-engager-pour-les-entreprises/modeles-economiques-soutenable-comment-placer-l-impact-positif-au-caeur-de-l-entreprise>
47. Gérer les conflits dans les DAOs : Mécanismes et outils pour une harmonie décentralisée, dernier accès : juillet 31, 2025, <https://w3r.one/fr/blog/blockchain-web3/daos/mecanisme-gouvernance/gestion-conflits-daos-mecanismes-outils-harmonie-decentralisee>
48. Decentralised autonomous organisation - Internet Policy Review, dernier accès : juillet 31, 2025, <https://policyreview.info/open-abstracts/decentralised-autonomous-organisation>
49. Position: Towards a Responsible LLM-empowered Multi-Agent Systems - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/html/2502.01714v1>
50. [2505.02077] Open Challenges in Multi-Agent Security: Towards Secure Systems of Interacting AI Agents - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/abs/2505.02077>

51. ASI-ARCH and the Double-Edged Sword of Self-Improving AI | The Neuron, dernier accès : juillet 31, 2025, https://www.theneuron.ai/explainer-articles/asi-arch-and-the-double-edged-sword-of-self-improving-ai?utm_source=www.theneurondaily.com&utm_medium=referral&utm_campaign=six-new-gpt-5-models-a-6k-robot-gymnast-and-an-ai-that-builds-ai
52. [2507.18074] AlphaGo Moment for Model Architecture Discovery - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/abs/2507.18074>
53. Qu'est-ce que le méta-apprentissage ? | IBM, dernier accès : juillet 31, 2025, <https://www.ibm.com/fr-fr/think/topics/meta-learning>
54. Méta-apprentissage : Comment les machines apprennent à apprendre - DataCamp, dernier accès : juillet 31, 2025, <https://www.datacamp.com/fr/blog/meta-learning>
55. [2505.08827] Self Rewarding Self Improving - arXiv, dernier accès : juillet 31, 2025, <https://arxiv.org/abs/2505.08827>
56. Computational Social Science Lab | Projects and Labs, dernier accès : juillet 31, 2025, <https://css.seas.upenn.edu/>
57. On agent-based modeling and computational social science - PMC, dernier accès : juillet 31, 2025, <https://pmc.ncbi.nlm.nih.gov/articles/PMC4094840/>
58. Computational Social Science and Sociology - PubMed, dernier accès : juillet 31, 2025, <https://pubmed.ncbi.nlm.nih.gov/34824489/>
59. La révolution des agents IA : 4 obstacles à surmonter pour les entreprises - Galaksiya.com, dernier accès : juillet 31, 2025, <https://www.galaksiya.com/fr/articles/the-ai-agent-revolution-4-barriers-enterprises-must-overcome>