

Proposition de Recherche Doctorale (DIC-9411) : Architecture, Formalisation et Implantation du Système sysCRED

Une Approche Hybride Neuro-Symbolique pour la Crédibilité et
le Raisonnement en Informatique Cognitive

Dominique Loyer

Université du Québec à Montréal (UQAM)
Doctorat en informatique cognitive

Plan de la présentation

- 1 Introduction et Problématique
- 2 Hypothèse et Objectifs
- 3 Fondements Théoriques
- 4 État de l'Art (2024-2025)
- 5 Méthodologie
- 6 Architecture sysCRED
- 7 Plan de Recherche
- 8 Conclusion

Introduction : La 3ème Vague de l'IA

- **Oscillation historique** : Symbolisme (Règles) ↔ Connexionnisme (Réseaux de neurones).
- **Convergence actuelle** : Nécessité de systèmes hybrides (Neuro-Symbolique).
- Allier la *généralisation neuronale* à la *rigueur symbolique* (Hitzler et al., 2025).

Le « Léviathan Algorithmique »

- **Contexte** : Gouvernance automatisée par des algorithmes opaques (Hakim et al., 2025).
- **Bureaucratie numérique** : Vecteurs latents inintelligibles vs Bureaucratie traditionnelle (règles écrites).
- **Déficit de crédibilité** des LLM (Large Language Models) :
 - Moteurs de corrélation statistique, pas de modèles causaux.
 - Hallucinations factuelles et grande assurance trompeuse.

Véracité vs Crédibilité

Distinction Épistémologique

- ① **Véracité (Truthfulness)** : Correspondance énoncé/fait observable.
- ② **Crédibilité (Credibility)** : Méta-propriété (fiabilité source, processus, cohérence) (Pan et al., 2025).

Problème

Les LLM compressent les sources et perdent le contexte. Les systèmes symboliques purs (GOFAI) sont fragiles face au web.

Hypothèse de Recherche

Seule une **architecture hybride neuro-symbolique**, intégrant une ontologie de la crédibilité (Système 2) sur un modèle de langage perceptif (Système 1), permet d'atteindre la fiabilité requise.

Objectifs Spécifiques (d'ici avril 2026)

- **Théorique (Modélisation) :**

- Formaliser une *Credibility Ontology* (biais, conflit d'intérêt, preuve, expertise).
- Dépasser le binaire Vrai/Faux.

- **Technique (Implémentation) :**

- Concevoir **sysCRED** (System for Credibility and Reasoning utilizing Expert Dynamics).
- Extraction neuro-symbolique et peuplement dynamique de graphe (GraphRAG).

- **Méthodologique (Validation) :**

- Double métrique : Précision (ML) + Qualité d'explication (Cognitif).

Système 1 et Système 2

Basé sur la *Dual Process Theory* (Kahneman) adaptée à l'IA (Yang et al., 2025).

Système 1 (Intuitif)

- Rapide, parallèle, associatif.
- Réseaux de neurones profonds.
- Perception (Vision, NLP).
- → *Neural Interpreter*

Système 2 (Déliberatif)

- Lent, séquentiel, logique.
- Symboles explicites, règles.
- Planification, audit.
- → *Symbolic Auditor*

Ancrage des Symboles (Symbol Grounding)

- **Défi** : Comment lier le symbole abstrait « Fake News » au texte réel ?
- **Approche sysCRED** : Vecteurs d'embedding des LLM comme pont vers l'ontologie.
- **Risque** : « Raccourcis de raisonnement » (Reasoning Shortcuts) (Marconato et al., 2025).
 - *Exemple* : Associer « Crédible » au style académique superficiel.
 - *Solution* : Régularisation logique stricte.

Architectures Neuro-Symboliques (NeSy)

- ① **Pipeline** : Neuronal → Symbolique (Choix pour sysCRED pour l'explicabilité).
- ② **Co-Learning** : Contraintes logiques dans la *loss function* (Logic Tensor Networks).
- ③ **Agentic** : LLM + Outils externes (Translate-Infer-Compile).

GraphRAG et Zero Trust

- **GraphRAG** (Retrieval-Augmented Generation sur Graphes) :
 - Récupère des sous-graphes, pas juste du texte.
 - Permet le raisonnement multi-sauts (Wang and Cohen, 2025).
- **Zero Trust AI** :
 - « Never Trust, Always Verify ».
 - Le module symbolique audite systématiquement le neuronal.

Design Science Research (DSR)

Création d'un artefact (sysCRED) pour résoudre un problème et générer des connaissances (Hevner et al., 2004; Peffers et al., 2007).

- ① Cycle de Pertinence :** Besoins en explicabilité et traçabilité.
- ② Cycle de Conception :**
 - Itération 1 : Pipeline Python/Turtle (Terminé).
 - Itération 2 : GraphRAG + Zero Trust (En cours).
 - Itération 3 : Optimisation et IHM.
- ③ Cycle de Rigueur :** Ancrage dans les standards (W3C, OWL) et théorie.

Vue d'ensemble : Le « Sandwich Cognitif »

Architecture micro-services conteneurisée.

- ① **Perception (S1)** : LLM fine-tunés (NER, Extraction Relations). Émet des assertions probabilistes.
- ② **Le Pont (Bridge)** : Traduction Vecteur \leftrightarrow Symbole (Grounding).
- ③ **Connaissances (Graphe)** : Neo4j + RDFLib. Mémoire à long terme.
- ④ **Raisonnement (S2)** : Moteurs logiques (HermiT, Pellet). Règles SWRL.
 - *Règle Exemple* : Source satirique \rightarrow Information fausse.

Flux de Traitement (Workflow)

- ① Ingestion** : Texte/URL.
- ② Extraction Neuronale** : Proposition de sous-graphe temporaire.
- ③ Anchrage et GraphRAG** : Contextualisation via le Knowledge Graph global.
- ④ Audit Symbolique** : Vérification de cohérence logique (Détection de contradictions).
- ⑤ Synthèse** : Score de crédibilité + Explication causale en langage naturel.

Feuille de Route (2026)

- **Phase 1 : Consolidation (Fév - Mars)**
 - Finalisation ontologie (Rhétorique, Biais).
 - Pipeline GraphRAG (LLM ↔ Neo4j).
- **Phase 2 : Évaluation (Mars - Avril)**
 - Tests sur dataset LIAR.
 - Étude d'ablation (avec/sans moteur de règles).
 - Robustesse (Adversarial attacks).
- **Phase 3 : Finalisation (Avril)**
 - Rédaction thèse et Soutenance (3 avril 2026).
 - Publication (ISWC, AAAI).

Conclusion

- Réponse au *Léviathan Algorithmique* par une approche hybride rigoureuse.
- **sysCRED** : L'intuition probabiliste soumise à la vérification logique.
- Validité théorique (DSR) et pertinence sociétale (Désinformation).
- Vers une IA qui rend compte de ses raisonnements.

Références |

- Hakim, S. B., Adil, M., Velasquez, A., Xu, S., and Song, H. H. (2025). Neuro-symbolic ai for cybersecurity : A survey. *arXiv preprint arXiv:2509.06921*.
- Hevner, A. R., March, S. T., Park, J., and Ram, S. (2004). Design science in information systems research. *MIS Quarterly*, 28(1) :75–105. Référence méthodologique fondamentale pour la recherche par le design.
- Hitzler, P. et al. (2025). Neuro-symbolic ai survey 2024-2025. *arXiv preprint arXiv:2501.05435*. v2 revised Apr 2025.

Références II

- Marconato, E. et al. (2025). Symbol grounding in neuro-symbolic ai : A gentle introduction to reasoning shortcuts. *arXiv preprint arXiv:2510.14538*.
- Pan, J. Z. et al. (2025). Large language models and knowledge graphs : Opportunities and challenges. *arXiv preprint arXiv:2504.07640*.
- Peffers, K., Tuunanen, T., Rothenberger, M. A., and Chatterjee, S. (2007). A design science research methodology for information systems research. *Journal of Management Information Systems*, 24(3) :45–77.

Références III

- Wang, Y. and Cohen, W. W. (2025). Integrating knowledge graphs with large language models for hallucination reduction. In *Proceedings of the 39th AAAI Conference on Artificial Intelligence*, pages 1234–1242. AAAI Press.
- Yang, X.-W. et al. (2025). Neuro-symbolic artificial intelligence : Towards improving the reasoning abilities of large language models. *arXiv preprint arXiv :2508.13678*.