

Analyse Architecturale DeepSeek

Modèle 4C Complet

Chaabane Ghazlene & Mhiri Nour & Morghom Nada & Sahnouni Rahma

Date : 26 septembre 2025

Version : 1.0

Table des matières

1	Introduction	3
2	Diagramme de Contexte (C1)	3
2.1	Description Détaillée	3
2.1.1	Éléments Clés	3
2.1.2	Spécificités DeepSeek	3
2.2	Diagramme	3
3	Diagramme de Conteneurs (C2)	4
3.1	Description Détaillée	4
3.1.1	Architecture Microservices	4
3.1.2	Avantages Architecturaux	4
3.2	Diagramme	4
4	Diagramme de Composants (C3)	6
4.1	Description Détaillée	6
4.1.1	Composants Principaux	6
4.2	Diagramme	6
5	Diagramme de Code (C4)	7
5.1	Description Détaillée	7
5.1.1	Modules Principaux	7
5.2	Diagramme	7
6	Analyse Technique Détaillée	9
6.1	Moteur d'Inférence Statistique	9
6.1.1	Fonctionnement Probabiliste	9
6.1.2	Optimisations Matérielles	9
6.2	Gestion des Cas Limites	9
6.2.1	Questions Inconnues	9
6.2.2	Détection d'Hallucinations	9

7	Synthèse Architecturale	9
7.1	Points Forts	9
7.2	Intégration des Concepts	10
7.3	Perspectives d'Évolution	10
8	Conclusion	10

1 Introduction

Cette analyse présente l'architecture de DeepSeek selon le modèle 4C (Context, Containers, Components, Code), enrichie des concepts techniques discutés lors de notre échange sur le fonctionnement interne du modèle.

2 Diagramme de Contexte (C1)

2.1 Description Détaillée

Objectif : Vue d'ensemble de l'écosystème DeepSeek et de ses interactions externes.

2.1.1 Éléments Clés

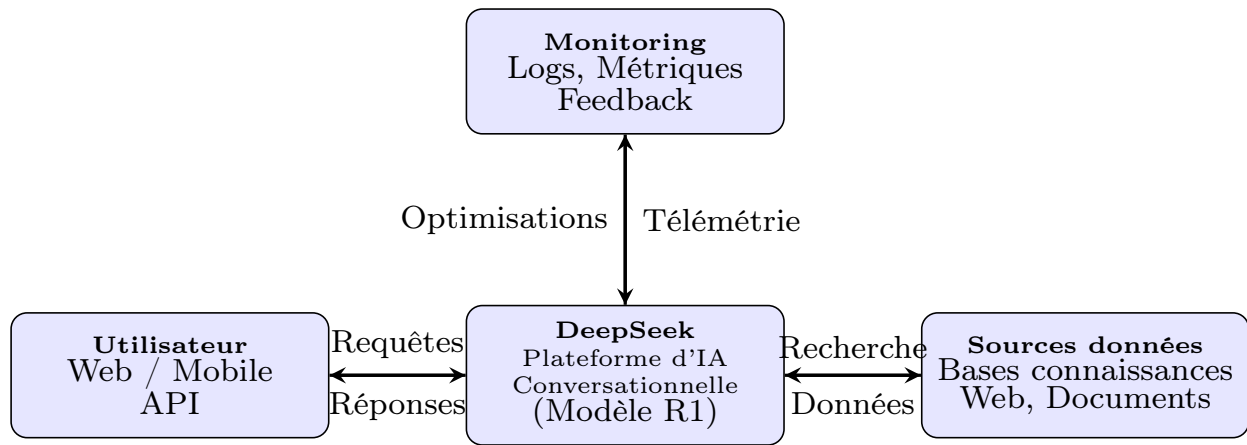
- **DeepSeek comme boîte noire :** Représenté comme une unité fonctionnelle unique
- **Acteurs externes :** Utilisateurs, sources de données, systèmes de monitoring

2.1.2 Spécificités DeepSeek

- Support multi-modal (web, mobile, API)
- Intégration de recherches externes pour actualisation
- Monitoring actif pour adaptation continue

2.2 Diagramme

Ce diagramme présente l'architecture contextuelle de DeepSeek, une plateforme d'intelligence artificielle conversationnelle basée sur le modèle R1. Au centre se trouve le système DeepSeek qui orchestre l'ensemble des interactions avec trois composants externes clés. Les utilisateurs accèdent à la plateforme via différents canaux (interfaces web, applications mobiles ou API) pour soumettre leurs requêtes conversationnelles et recevoir les réponses générées par l'IA. Le système interroge dynamiquement diverses sources de données externes comprenant des bases de connaissances, des contenus web et des documents pour enrichir ses réponses avec des informations pertinentes et à jour. Un système de monitoring supervise l'ensemble en collectant la télémétrie opérationnelle (logs, métriques de performance et feedback utilisateur) depuis DeepSeek, puis retourne des optimisations pour améliorer continuellement les performances et la qualité des réponses. Cette architecture illustre un écosystème d'IA conversationnelle moderne où l'apprentissage continu et l'amélioration des performances sont assurés par une boucle de feedback permanent entre les interactions utilisateur, l'accès aux données externes et le monitoring intelligent.



3 Diagramme de Conteneurs (C2)

3.1 Description Détaillée

Objectif : Découpage en services déployables indépendamment.

3.1.1 Architecture Microservices

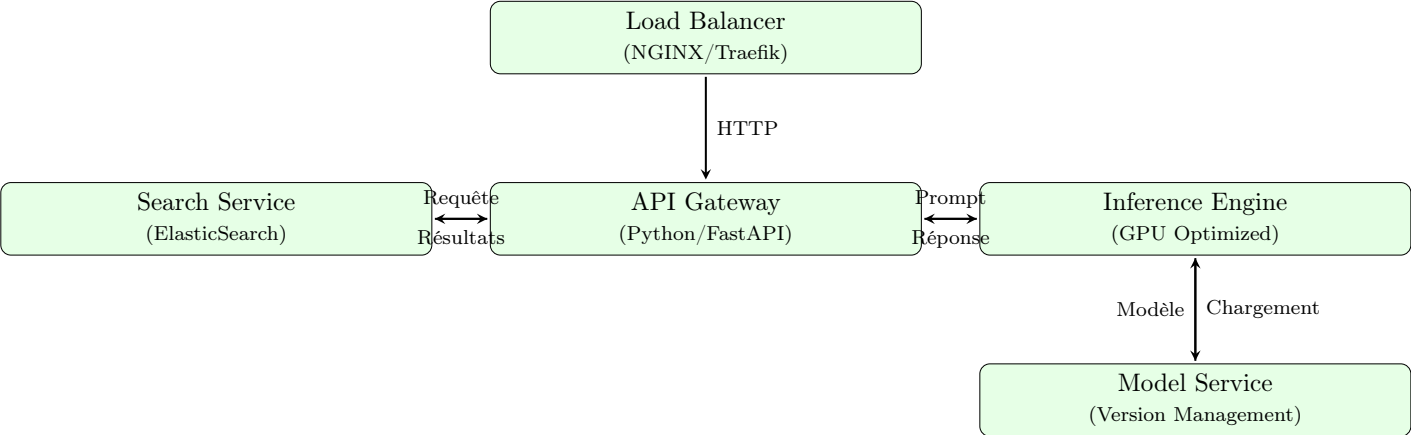
- **Load Balancer** : Gestion du trafic et haute disponibilité
- **API Gateway** : Point d'entrée unique avec sécurisation
- **Inference Engine** : Cœur de l'IA sur infrastructure GPU

3.1.2 Avantages Architecturaux

- **Scalabilité Horizontale** : Chaque service peut scale indépendamment
- **Résilience** : Isolation des fautes entre composants

3.2 Diagramme

Ce diagramme illustre l'architecture en conteneurs de DeepSeek, montrant comment les différents services collaborent pour traiter les requêtes utilisateur. Le flux commence par le Load Balancer (NGINX/Traefik) qui distribue le trafic HTTP entrant vers l'API Gateway central (Python/FastAPI), point d'orchestration de toutes les interactions. Lorsqu'une requête arrive, l'API Gateway coordonne deux services spécialisés : il interroge le Search Service (ElasticSearch) pour récupérer des informations contextuelles pertinentes dans les bases de données, puis transmet le prompt enrichi à l'Inference Engine optimisé GPU pour la génération de réponses. L'Inference Engine s'appuie sur le Model Service pour charger dynamiquement les versions appropriées des modèles d'IA selon les besoins, permettant une gestion flexible des différentes variantes et mises à jour de modèles. Cette architecture modulaire assure une séparation claire des responsabilités : la recherche d'information, l'inférence IA et la gestion des modèles sont découplées, permettant une scalabilité indépendante de chaque composant et une maintenance optimisée du système global.



4 Diagramme de Composants (C3)

4.1 Description Détaillée

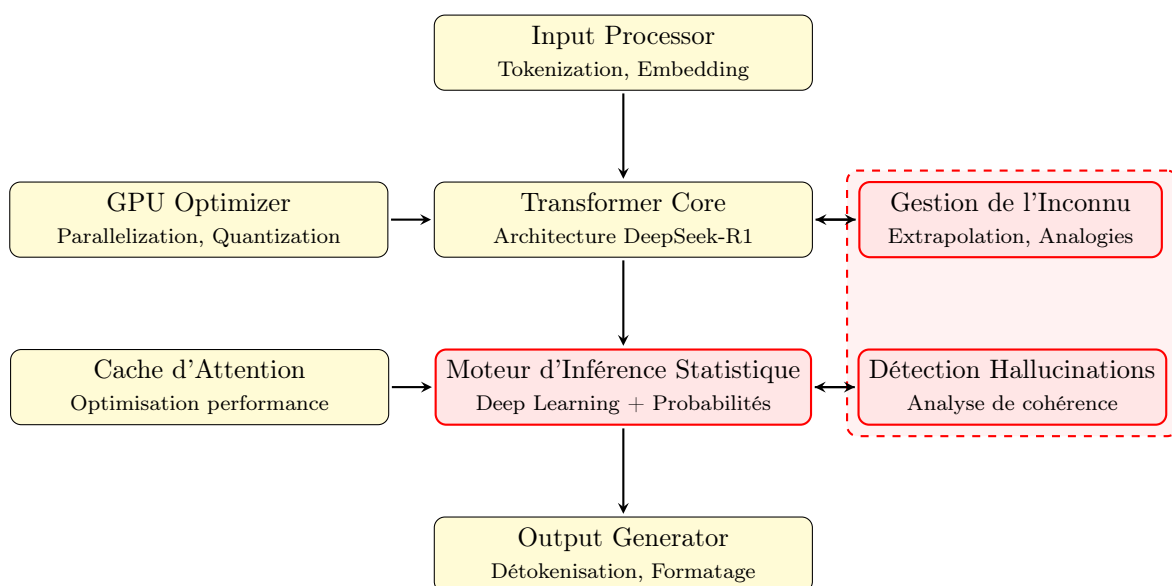
Objectif : Architecture interne du moteur d'inférence DeepSeek.

4.1.1 Composants Principaux

- **Transformer Core** : Architecture DeepSeek-R1
- **Statistical Inference Engine** : Génération probabiliste
- **Gestion des cas limites** : Extrapolation et détection d'hallucinations

4.2 Diagramme

Ce diagramme détaille l'architecture interne des composants de DeepSeek, révélant le pipeline de traitement sophistiqué du modèle R1. Le flux principal suit une séquence linéaire : l'Input Processor effectue la tokenisation et la génération d'embeddings des requêtes utilisateur, puis le Transformer Core (architecture DeepSeek-R1) traite ces représentations vectorielles pour comprendre le contexte et les relations sémantiques. Le Moteur d'Inférence Statistique combine ensuite les techniques de deep learning avec des calculs probabilistes pour générer les prédictions, avant que l'Output Generator ne reconvertisse les tokens en texte formaté pour l'utilisateur final. Deux composants spécialisés (encadrés en rouge) gèrent les défis critiques de l'IA générative : le module de Gestion de l'Inconnu utilise l'extrapolation et les analogies pour traiter les concepts non vus pendant l'entraînement, tandis que la Détection d'Hallucinations analyse la cohérence des réponses pour minimiser les erreurs factuelles. Des optimisations techniques accompagnent ce processus : le GPU Optimizer assure la parallélisation et la quantisation pour des performances maximales, et le Cache d'Attention stocke les patterns fréquents pour accélérer l'inférence, créant ainsi un système à la fois puissant et fiable.



5 Diagramme de Code (C4)

5.1 Description Détaillée

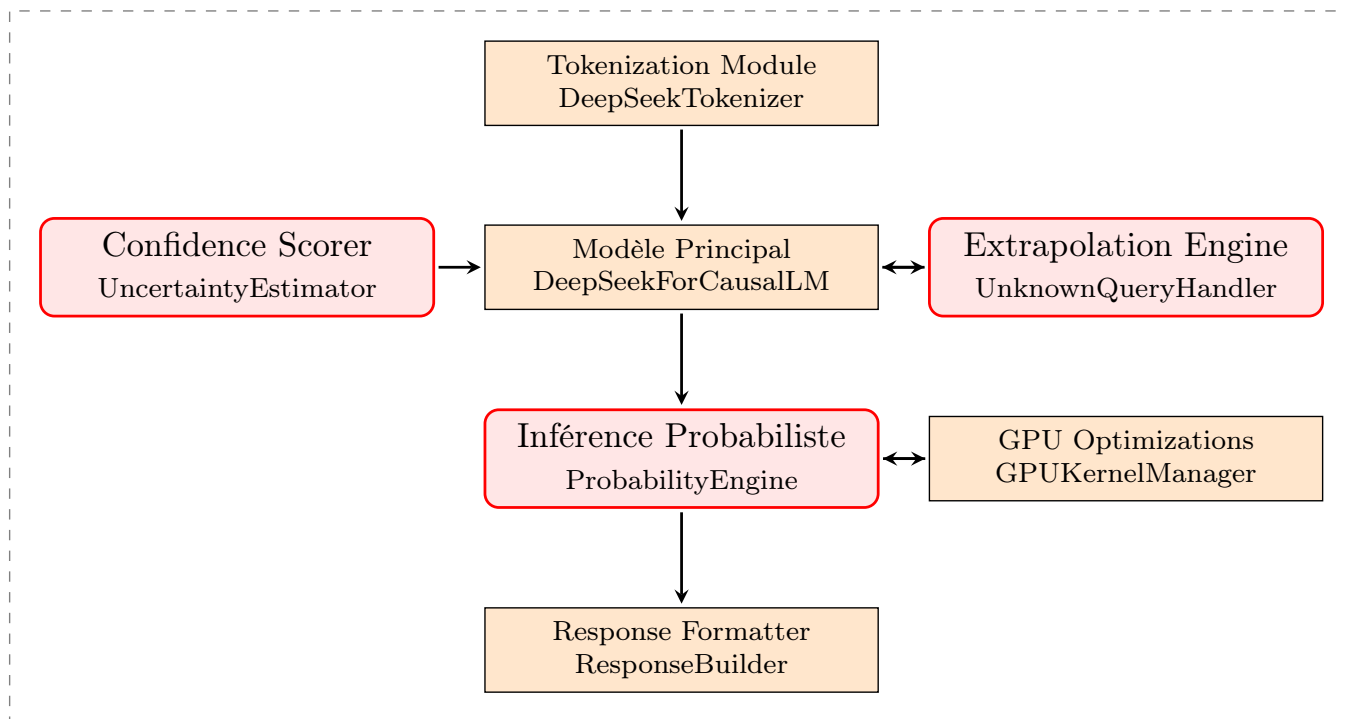
Objectif : Implémentation détaillée des modules logiciels.

5.1.1 Modules Principaux

- **DeepSeekTokenizer** : Gestion du vocabulaire et encodage
- **ProbabilityEngine** : Implémentation sampling et beam search
- **GPU Optimizations** : Kernels CUDA/ROCM optimisés

5.2 Diagramme

Ce diagramme présente l'implémentation détaillée au niveau code des modules DeepSeek, révélant l'architecture logicielle concrète du système d'IA conversationnelle. Le pipeline central suit un flux séquentiel strict : le DeepSeekTokenizer convertit les entrées textuelles en tokens numériques, puis le modèle principal DeepSeekForCausalLM (implémentation du transformer R1) traite ces séquences pour générer les représentations contextuelles. Le ProbabilityEngine effectue ensuite l'inférence probabiliste en calculant les distributions de probabilités sur le vocabulaire, avant que le ResponseBuilder ne formate la réponse finale pour l'utilisateur. Trois modules spécialisés enrichissent ce processus : l'UnknownQueryHandler gère l'extrapolation pour les requêtes hors domaine d'entraînement, l'UncertaintyEstimator évalue le niveau de confiance des prédictions pour identifier les réponses potentiellement incertaines, et le GPUIKernelManager optimise l'exécution sur processeurs graphiques pour maximiser les performances. Cette architecture modulaire (encadrée en pointillés) assure une séparation claire entre la logique métier, les optimisations techniques et la gestion de l'incertitude, permettant une maintenance et une évolutivité optimales du code source.



6 Analyse Technique Détaillée

6.1 Moteur d'Inférence Statistique

Le cœur de l'architecture DeepSeek repose sur un moteur d'inférence statistique qui diffère des systèmes symboliques traditionnels.

6.1.1 Fonctionnement Probabiliste

- **Inférence basée sur les probabilités** : Calcul du token suivant le plus probable
- **Sampling et Beam Search** : Algorithmes de génération de texte
- **Attention Mechanism** : Calcul des relations contextuelles

6.1.2 Optimisations Matérielles

- **Parallelization GPU** : Calculs massivement parallèles
- **Quantification** : Réduction de précision pour la performance
- **Memory Optimization** : Gestion efficace de la mémoire

6.2 Gestion des Cas Limites

6.2.1 Questions Inconnues

- **Extrapolation** : Génération basée sur des concepts similaires
- **Raisonnement Analogique** : Utilisation de patterns appris
- **Estimation d'Incertitude** : Calcul du niveau de confiance

6.2.2 Détection d'Hallucinations

- **Analyse de Cohérence** : Vérification interne des réponses
- **Validation Contextuelle** : Comparaison avec le contexte
- **Scoring de Confiance** : Attribution d'un score de fiabilité

7 Synthèse Architecturale

7.1 Points Forts

- **Performance** : Optimisations GPU avancées
- **Scalabilité** : Architecture microservices
- **Robustesse** : Gestion des erreurs intégrée
- **Évolutivité** : Gestion des versions de modèles

7.2 Intégration des Concepts

- **Inférence Statistique** : Fondement du raisonnement
- **Deep Learning** : Apprentissage profond intégral
- **Optimisations** : Performance et efficacité

7.3 Perspectives d'Évolution

- **Multi-modalité** : Intégration image/audio
- **Apprentissage Continu** : Adaptation en temps réel
- **Optimisations** : Nouvelles architectures matérielles

8 Conclusion

L'architecture DeepSeek démontre une conception moderne qui intègre efficacement les avancées récentes en deep learning avec les bonnes pratiques d'ingénierie logicielle. La combinaison d'une infrastructure scalable, d'optimisations matérielles avancées et de mécanismes de gestion robustes des cas limites positionne cette plateforme comme une solution d'IA conversationnelle de haute qualité.