

## Building LLM Applications with Prompt Engineering

Durée: 1 Jour    Réf de cours: GK847008    Méthodes d'apprentissage: Intra-entreprise & sur-mesure

---

### Résumé:

Very large deep neural networks (DNNs), whether applied to natural language processing (e.g., GPT-3), computer vision (e.g., huge Vision Transformers), or speech AI (e.g., Wave2Vec 2) have certain properties that set them apart from their smaller counterparts. As DNNs become larger and are trained on progressively larger datasets, they can adapt to new tasks with just a handful of training examples, accelerating the route toward general artificial intelligence. Training models that contain tens to hundreds of billions of parameters on vast datasets isn't trivial and requires a unique combination of AI, high-performance computing (HPC), and systems knowledge.

### Company Events

These events can be delivered exclusively for your company at our locations or yours, specifically for your delegates and your needs. The Company Events can be tailored or standard course deliveries.

---

### Objectifs pédagogiques:

- Train neural networks across multiple servers
  - Use techniques such as activation checkpointing, gradient accumulation, and various forms of model parallelism to overcome the challenges associated with large-model memory footprint
  - Capture and understand training performance characteristics to optimize model architecture
  - Deploy very large multi-GPU models to production using NVIDIA Triton Inference Server
-

## Contenu:

### Module 1: Introduction

- Orienter vers les principaux sujets du atelier, l'emploi du temps et les prérequis.
- Apprendre pourquoi l'ingénierie des prompts est au cœur de l'interaction avec les Modèles de Langage de Grande Taille (LLMs).
- Discuter de la façon dont l'ingénierie des prompts peut être utilisée pour développer de nombreuses classes d'applications basées sur les LLM.
- Apprendre sur NVIDIA LLM NIM, utilisé pour déployer le Llama 3.1 LLM utilisé dans l'atelier.

### Module 2: Introduction to Prompting

- Devenir familier avec l'environnement de l'atelier.
- Créer et visualiser des réponses à partir de vos premiers prompts en utilisant l'API OpenAI et LangChain.
- Apprendre comment streamer les réponses des LLM, et envoyer des prompts aux LLMs par lots, en comparant les différences de performance.
- Commencer à pratiquer le processus de développement itératif des prompts.
- Créer et utiliser vos premiers modèles de prompts.
- Faire un mini-projet où vous effectuez une combinaison d'analyse et de tâches génératives sur un lot d'entrées.

### Module 3: LangChain Expression Language (LCEL), Runnable, et Chains

- Apprendre sur les Runnable de LangChain, et la capacité de les composer en chaînes en utilisant le LangChain Expression Language (LCEL).
- Écrire des fonctions personnalisées et les convertir en Runnable qui peuvent être incluses dans les chaînes de LangChain.
- Composer plusieurs chaînes LCEL en une seule chaîne d'application plus grande.
- Exploiter les opportunités de travail parallèle en composant des chaînes LCEL parallèles.
- Faire un mini-projet où vous effectuez une combinaison d'analyse et de tâches génératives sur un lot d'entrées en utilisant LCEL et l'exécution parallèle.

### Module 4: Prompting With Messages

- Apprendre sur deux des principaux types de messages, humains et IA, et comment les utiliser explicitement dans le code d'application.
- Fournir des modèles de chat avec des exemples instructifs en utilisant une technique appelée prompting à tirage limité.
- Travailler explicitement avec le message système, qui vous permettra de définir un personnage et un rôle globaux pour vos modèles de chat.
- Utiliser le prompting à chaîne de pensée pour améliorer la capacité des LLMs à effectuer des tâches nécessitant un raisonnement complexe.
- Gérer les messages pour conserver l'historique de la conversation et activer la fonctionnalité de chatbot.
- Faire un mini-projet où vous construisez une application de chatbot simple mais flexible capable d'assumer divers rôles.

### Module 5: Structured Output

- Explorer quelques méthodes de base pour utiliser les LLMs afin de générer des données structurées en lot pour une utilisation en aval.
- Générer une sortie structurée en combinant des classes Pydantic et LangChain's `JsonOutputParser`.
- Apprendre comment extraire des données et les tagger à partir de textes longs.
- Faire un mini-projet où vous utilisez des techniques de génération de données structurées pour effectuer l'extraction de données et l'étiquetage de documents non structurés.

### Module 6: Tool Use and Agents

- Créer une fonctionnalité externe aux LLMs appelée outils, et rendre les LLMs conscients de leur disponibilité.
- Créer un agent capable de raisonner sur l'utilisation appropriée des outils, et intégrer les résultats de l'utilisation des outils dans ses réponses.
- Faire un mini-projet où vous créez un agent LLM capable d'utiliser des appels API externes pour améliorer ses réponses avec des données en temps réel.

### Module 7: Assessment and Final Review

- Réviser les apprentissages clés et répondre aux questions.
- Gagner un certificat de compétence pour l'atelier.
- Compléter le sondage de l'atelier.
- Obtenir des recommandations pour les prochaines étapes de votre parcours d'apprentissage.

## Autres moyens pédagogiques et de suivi:

- Compétence du formateur : Les experts qui animent la formation sont des spécialistes des matières abordées et ont au minimum cinq ans d'expérience d'animation. Nos équipes ont validé à la fois leurs connaissances techniques (certifications le cas échéant) ainsi que leur compétence pédagogique.
- Suivi d'exécution : Une feuille d'embarquement par demi-journée de présence est signée par tous les participants et le formateur.
- En fin de formation, le participant est invité à s'auto-évaluer sur l'atteinte des objectifs énoncés, et à répondre à un questionnaire de satisfaction qui sera ensuite étudié par nos équipes pédagogiques en vue de maintenir et d'améliorer la qualité de nos prestations.

Délais d'inscription :

- Vous pouvez vous inscrire sur l'une de nos sessions planifiées en inter-entreprises jusqu'à 5 jours ouvrés avant le début de la formation