

Agents IA à base de LLM avec Langchain

Christophe Guyeux

IUT NFC – UMLP

10 mai 2026

- 1 Introduction
- 2 Premier prompt avec un LLM
- 3 Sorties structurées avec Pydantic
- 4 Outils et appels de fonctions
- 5 Premier agent avec `create_agent`
- 6 Mémoire et streaming
- 7 LangGraph pour les boucles complexes
- 8 RAG minimal

Introduction

Objectifs du cours

- Comprendre ce qu'est un **LLM** et un **agent** à base de LLM.
- Maîtriser les **briques de Langchain** : modèles, prompts, chaînes, outils, agents.
- Construire **progressivement** un agent capable de raisonner et d'utiliser des outils.
- **Tout coder** avec des modèles *gratuits* : Mistral (cloud, avec clé) ou Ollama (local).

Qu'est-ce qu'un LLM ?

Un **Large Language Model** est un réseau de neurones (Transformer) entraîné à prédire le mot suivant à partir d'un contexte.

- Entrée : une séquence de tokens (texte découpé).
- Sortie : une distribution de probabilités sur le prochain token.
- Itéré, cela produit une **génération** cohérente.

Conséquence pratique

Un LLM ne « sait » rien : il *prédit du texte plausible*. D'où l'intérêt de lui donner du contexte (prompts, outils, mémoire) — c'est tout l'objet du cours.

Qu'est-ce qu'un agent ?

Un **agent** est un programme qui :

- 1 reçoit un objectif en langage naturel,
- 2 décide d'actions à entreprendre (appeler un outil, consulter une base, écrire un fichier...),
- 3 observe les résultats et **itère** jusqu'à atteindre l'objectif.

Le LLM joue le rôle de **cerveau décisionnel** ; les outils sont ses « mains ».

Objectif → LLM raisonne → outil → observation → LLM raisonne → ...

Pourquoi Langchain ?

Langchain est un framework Python qui standardise l'intégration des LLMs :

- **Agnostique** : Mistral, OpenAI, Anthropic, DeepSeek, Ollama, etc. derrière une interface unique.
- **Composable** (LCEL) : `prompt | model | parser` comme dans un shell Unix.
- **Outillage agent** : `create_agent`, gestion de la mémoire, des outils, du streaming.
- **Écosystème** : `langgraph` (boucles complexes), `langchain-community`, RAG, vecteurs.

Référence officielle : <https://docs.langchain.com/oss/python/langchain/agents>

- 1 **Premier prompt + LLM** : invocation simple, prompt template, chaîne LCEL.
- 2 Sorties structurées avec Pydantic.
- 3 Outils (*tools*) et appels de fonctions.
- 4 Premier agent avec `create_agent`.
- 5 Mémoire et streaming.
- 6 LangGraph pour les boucles agentiques complexes.
- 7 RAG minimal : retrouver l'information dans des documents.

Prérequis et installation

- Python **3.11+**
- uv (<https://docs.astral.sh/uv/>)
- Optionnel : Ollama (<https://ollama.com/>) pour exécuter un LLM en local

```
# Installation initiale (chapitre 00)
```

```
uv sync
```

Premier prompt avec un LLM

Au début de chaque cours, on rappelle les paquets à installer :

```
# Bibliothèques à installer (chapitre 01)

uv add langchain langchain-mistralai python-dotenv

# Variante locale sans clé API (Ollama) :
uv add langchain-ollama
ollama pull qwen2.5:3b
```

Configurer la clé d'API

Créez un fichier `.env` à la racine du projet :

```
1 LLM_BACKEND=mistral
2 MISTRAL_API_KEY=xxxxxxxxxxxxxxxxxxxxxxxxxxxx
```

- Mistral : créez une clé sur <https://console.mistral.ai/> (tier *Experiment* gratuit, sans CB).
- **Tout en local** : passez à Ollama (`LLM_BACKEND=ollama`) après `ollama pull qwen2.5:3b`.

La bibliothèque `python-dotenv` chargera ces variables au lancement du script.

Première invocation d'un LLM

Le pattern minimal : on instancie un modèle, on l'**invoque** avec une chaîne de caractères, on lit `response.content`.

```
1  """Première invocation : un prompt, un LLM, une réponse."""
2  import os
3  from dotenv import load_dotenv
4
5  load_dotenv()
6  backend = os.getenv("LLM_BACKEND", "mistral")
7
8  if backend == "mistral":
9      from langchain_mistralai import ChatMistralAI
10     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
11 else:
12     from langchain_ollama import ChatOllama
13     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
14
15 response = llm.invoke("Quelle est la capitale de l'Albanie ? Réponds en une phrase.")
16 print(response.content)
```

Que se passe-t-il ?

- 1 `load_dotenv()` : importe la clé d'API depuis `.env`.
- 2 `ChatMistralAI(...)` : instancie un *chat model* Langchain.
- 3 `llm.invoke(...)` : envoie un appel HTTP à Mistral, attend la réponse.
- 4 `response.content` : le texte généré par le modèle.

À retenir

`ChatXxx` est l'interface uniforme : changer de fournisseur ne demande **qu'une ligne**. Le reste du code reste identique.

Structurer le prompt : message système + utilisateur

Un ChatPromptTemplate sépare le **rôle** (système, utilisateur) du **contenu**, et permet l'interpolation de variables.

```
1 """ChatPromptTemplate : séparer rôle système et message utilisateur."""
2 import os
3 from dotenv import load_dotenv
4 from langchain_core.prompts import ChatPromptTemplate
5
6 load_dotenv()
7 backend = os.getenv("LLM_BACKEND", "mistral")
8
9 if backend == "mistral":
10     from langchain_mistralai import ChatMistralAI
11     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
12 else:
13     from langchain_ollama import ChatOllama
14     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
15
16 prompt = ChatPromptTemplate.from_messages([
17     ("system", "Tu es un expert en géographie. Tu réponds en une seule phrase concise."),
18     ("human", "{question}"),
19 ])
20
21 messages = prompt.invoke({"question": "Quelle est la capitale du Bhoutan ?"})
22 response = llm.invoke(messages)
23 print(response.content)
```

Pourquoi un message système ?

- Il **cadre** le comportement du modèle (ton, format, contraintes).
- Il est *persistant* dans la conversation, contrairement à un préfixe ajouté à la question.
- Il facilite les **tests** : changer le système sans toucher aux questions utilisateur.

Bonne pratique : un message système *court* et *spécifique* produit de meilleurs résultats qu'un long préambule.

Composer une chaîne avec LCEL

Le **Langchain Expression Language** compose des étapes avec l'opérateur | (pipe).

```
1 """LCEL : composer prompt | llm | parser avec l'opérateur pipe."""
2 import os
3 from dotenv import load_dotenv
4 from langchain_core.prompts import ChatPromptTemplate
5 from langchain_core.output_parsers import StrOutputParser
6
7 load_dotenv()
8 backend = os.getenv("LLM_BACKEND", "mistral")
9
10 if backend == "mistral":
11     from langchain_mistralai import ChatMistralAI
12     llm = ChatMistralAI(model="mistral-small-latest", temperature=0.7)
13 else:
14     from langchain_ollama import ChatOllama
15     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0.7)
16
17 prompt = ChatPromptTemplate.from_template(
18     "Raconte-moi une blague courte sur le sujet : {sujet}"
19 )
20
21 chain = prompt | llm | StrOutputParser()
22
23 for sujet in ["les pompiers", "les data scientists"]:
24     print(f"--- {sujet} ---")
25     print(chain.invoke({"sujet": sujet}))
```

```
prompt | llm | StrOutputParser()
```

- `prompt` : transforme un dict en messages.
- `llm` : transforme des messages en `AIMessage`.
- `StrOutputParser()` : extrait `.content` en `str`.

Chaque maillon est un **Runnable**. On peut les substituer, les paralléliser (`RunnableParallel`), les router (`RunnableBranch`).

Et si je veux un autre fournisseur ?

Langchain expose la même interface ChatModel pour tous les fournisseurs : **seul l'import et l'instanciation changent**. Quelques exemples :

OpenAI

uv add langchain-openai

```
# Clé : platform.openai.com/api-keys
import os
from langchain_openai import ChatOpenAI

llm = ChatOpenAI(
    model="gpt-4.1-mini",
    temperature=0,
    api_key=os.environ["OPENAI_API_KEY"],
)
print(llm.invoke("Bonjour !").content)
```

DeepSeek

uv add langchain-deepseek

```
# Clé : platform.deepseek.com
import os
from langchain_deepseek import
    ChatDeepSeek

llm = ChatDeepSeek(
    model="deepseek-chat",
    temperature=0,
    api_key=os.environ["DEEPSEEK_API_KEY"]
)
print(llm.invoke("Bonjour !").content)
```

Google Gemini

uv add langchain-google-genai

```
# Clé : aistudio.google.com/app/apikey
import os
from langchain_google_genai import
    ChatGoogleGenerativeAI

llm = ChatGoogleGenerativeAI(
    model="gemini-2.0-flash",
    temperature=0,
    google_api_key=os.environ["
        GOOGLE_API_KEY"],
)
print(llm.invoke("Bonjour !").content)
```

Le reste du code (prompt, chain, StrOutputParser, agents...) reste **identique** : c'est la promesse de l'interface unifiée.

- 1 Modifiez `01_invoke_simple.py` pour afficher également `response.usage_metadata` et observer le coût en tokens.
- 2 Dans `02_chat_prompt.py`, changez le message système pour produire une réponse **en anglais**, puis **en deux phrases**.
- 3 Dans `03_chain.py`, faites varier la température (0, 0.7, 1.5) et comparez la diversité des blagues.
- 4 Bascule en local : `LLM_BACKEND=ollama` et relancez les trois scripts. Que constatez-vous sur la qualité et la latence ?

Récapitulatif

- Un ChatModel s'**invoque** avec une string ou des messages.
- Un ChatPromptTemplate structure le prompt et permet l'interpolation.
- Le pipe | de LCEL compose prompt, modèle et parser en une seule **chaîne**.
- Changer de fournisseur (Mistral ↔ Ollama) ne change **qu'une ligne**.

Suite

Au prochain cours : forcer le LLM à produire des sorties **structurées** avec Pydantic.

Sorties structurées avec Pydantic

Bibliothèques pour ce chapitre

```
# Bibliothèques à installer (chapitre 02)  
  
# Aucune nouvelle dépendance : pydantic est fourni par langchain-core.  
# Tout est déjà disponible après 'uv sync'.
```

pydantic est déjà installé via langchain-core.

Pourquoi structurer la sortie ?

Un ChatModel renvoie par défaut du **texte libre** : pratique pour un humain, fragile pour un programme.

- Comment extraire *automatiquement* l'âge dans « Marie a 34 ans » ?
- Comment garantir qu'un classifieur renvoie bien "positif" | "négatif" | "neutre" et rien d'autre ?
- Comment connecter un LLM à une base de données ou une API qui attend du **JSON typé** ?

Idée

Décrire la **forme attendue** avec un modèle Pydantic, et laisser Langchain demander au LLM de produire un objet conforme.

Rappel express : Pydantic

`pydantic.BaseModel` est la brique standard de modélisation typée en Python : on déclare des classes avec annotations, et la validation est gratuite.

- Champs typés (`int`, `str`, `list[X]`, `Literal[...]`, modèles imbriqués...).
- `Field(description=...)` : la description sert de **prompt** pour le LLM.
- Contraintes : `ge`, `le`, `min_length`, `pattern`, etc.
- Validation à l'instanciation : un objet créé est garanti conforme.

Premier modèle structuré

```
1 """Sortie structurée : du texte vers un objet Pydantic typé."""
2 import os
3 from dotenv import load_dotenv
4 from pydantic import BaseModel, Field
5
6 load_dotenv()
7 backend = os.getenv("LLM_BACKEND", "mistral")
8
9 if backend == "mistral":
10     from langchain_mistralai import ChatMistralAI
11     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
12 else:
13     from langchain_ollama import ChatOllama
14     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
15
16
17 class Personne(BaseModel):
18     """Une personne identifiée dans un texte."""
19     nom: str = Field(description="Prénom et nom complet")
20     age: int = Field(description="Âge en années")
21
22
23 structured_llm = llm.with_structured_output(Personne)
24 resultat = structured_llm.invoke("Marie Dupont a 34 ans et habite à Paris.")
25
26 print(resultat)
27 print(f"Type : {type(resultat).__name__}")
28 print(f"Nom : {resultat.nom}")
29 print(f"Age : {resultat.age}")
```

Que se passe-t-il ?

- 1 `class Personne(BaseModel)` : on déclare la forme attendue.
- 2 `llm.with_structured_output(Personne)` : Langchain transforme le modèle en *schema*, l'envoi au LLM (par tool-calling ou JSON mode selon le fournisseur).
- 3 `.invoke(...)` renvoie directement une **instance de Personne**, déjà validée.

À retenir

Plus de `response.content` à parser à la main. `resultat.nom`, `resultat.age` sont **typés** et disponibles tout de suite.

Champs imbriqués et listes

Un modèle peut contenir des listes, des sous-modèles, des contraintes (ge, le, pattern...).

```
1 """Extraction d'un objet avec listes et contraintes : critique de film."""
2 import os
3 from dotenv import load_dotenv
4 from pydantic import BaseModel, Field
5
6 load_dotenv()
7 backend = os.getenv("LLM_BACKEND", "mistral")
8
9 if backend == "mistral":
10     from langchain_mistralai import ChatMistralAI
11     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
12 else:
13     from langchain_ollama import ChatOllama
14     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
15
16 class Critique(BaseModel):
17     titre: str = Field(description="Titre du film")
18     note: int = Field(description="Note sur 10", ge=0, le=10)
19     points_forts: list[str] = Field(description="Aspects positifs cités")
20
21
22 texte = "Inception : effets visuels sublimes, scénario captivant. Je mets 9/10."
23
24
25 structured_llm = llm.with_structured_output(Critique)
26 critique = structured_llm.invoke(f"Extrais une critique de :\n{texte}")
27
28 print(f"{critique.titre} : {critique.note}/10")
29 for p in critique.points_forts:
30     print(f" - {p}")
```

Anatomie de l'extraction

- `titre: str` et `note: int` : champs scalaires typés.
- `ge=0, le=10` : contraintes de plage. Le LLM voit le schéma et s'y conforme ; Pydantic vérifie à la réception.
- `points_forts: list[str]` : un champ liste, parfait pour les énumérations extraites du texte.

Bonne pratique

Les `description=...` de chaque `Field` se retrouvent **telles quelles** dans le prompt envoyé au modèle. Soigner ces descriptions = soigner le prompt.

Classification avec Literal

typing.Literal[...] contraint un champ à un **ensemble fini** de valeurs. Idéal pour la classification.

```
1 """Classification : Literal contraint le LLM à choisir parmi des valeurs fixes."""
2 import os
3 from typing import Literal
4 from dotenv import load_dotenv
5 from pydantic import BaseModel, Field
6
7 load_dotenv()
8 backend = os.getenv("LLM_BACKEND", "mistral")
9
10 if backend == "mistral":
11     from langchain_mistralai import ChatMistralAI
12     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
13 else:
14     from langchain_ollama import ChatOllama
15     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
16
17
18 class Sentiment(BaseModel):
19     polarite: Literal["positif", "negatif", "neutre"]
20     confiance: int = Field(description="Confiance 0-100", ge=0, le=100)
21
22
23 structured_llm = llm.with_structured_output(Sentiment)
24
25 for m in [
26     "Ce film m'a touché, je le recommande !",
27     "Service catastrophique, je ne reviendrai pas.",
28     "Le rapport est arrivé hier matin.",
29 ]:
30     s = structured_llm.invoke(m)
31     print(f"{m[:38]:38s} -> {s.polarite:8s} ({s.confiance}%)"
```

Sortie typée = sortie sûre

- Le LLM *ne peut pas* renvoyer une polarité hors de la liste : Pydantic refuserait l'objet.
- En pratique, le schéma envoyé au modèle inclut l'énumération : il choisit dans le menu, plutôt que de halluciner une valeur.
- Combiné à un `int` avec `ge/le`, on obtient un classifieur autodescriptif et typé en quelques lignes.

Alternative explicite : PydanticOutputParser

with_structured_output est magique : il choisit la bonne stratégie selon le fournisseur. Si on veut tout contrôler (instructions de format dans le prompt, modèles sans tool-calling...), on chaîne un PydanticOutputParser.

```
1 """Alternative explicite : PydanticOutputParser dans une chaîne LCEL."""
2 import os
3 from dotenv import load_dotenv
4 from pydantic import BaseModel, Field
5 from langchain_core.prompts import ChatPromptTemplate
6 from langchain_core.output_parsers import PydanticOutputParser
7
8 load_dotenv()
9 backend = os.getenv("LLM_BACKEND", "mistral")
10
11 if backend == "mistral":
12     from langchain_mistralai import ChatMistralAI
13     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
14 else:
15     from langchain_ollama import ChatOllama
16     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
17
18
19 class Recette(BaseModel):
20     plat: str = Field(description="Nom du plat")
21     ingredients: list[str] = Field(description="Ingrédients principaux")
22
23
24 parser = PydanticOutputParser(pydantic_object=Recette)
25 prompt = ChatPromptTemplate.from_messages([
26     ("system", "Tu réponds uniquement au format demandé."),
27     ("human", "{question}\n\n{format_instructions}"),
28 ])
29
30 chain = prompt | llm | parser
31 recette = chain.invoke({"question": "Recette d'une omelette aux fines herbes."})
32
33 print(f"Plat : {recette.plat}")
34 for i in recette.ingredients:
35     print(f" - {i}")
```

Quand utiliser quoi ?

- `with_structured_output` (Modèle) : choix par défaut, concis.
- `PydanticOutputParser` dans une **LCEL chain** : quand on veut visualiser/maîtriser les `format_instructions`, ou cibler un modèle qui n'expose pas de tool-calling.
- Pour les sorties très simples : `StrOutputParser`, `JsonOutputParser`, ou rien du tout.

- 1 Modifiez `01_pydantic_basics.py` pour ajouter un champ `ville: str | None`. Testez sur une phrase qui ne mentionne aucune ville.
- 2 Dans `02_extraction.py`, ajoutez un sous-modèle `Acteur` avec `nom` et `role`, et un champ `acteurs: list[Acteur]`.
- 3 Dans `03_classification.py`, élargissez `Literal` à cinq valeurs (`très_positif`, `positif`, `neutre`, `négatif`, `très_négatif`). Que constatez-vous sur la qualité?
- 4 Comparez les `format_instructions` générés par `PydanticOutputParser.get_format_instructions()` pour deux modèles différents.

Récapitulatif

- Un modèle Pydantic décrit la **forme attendue** de la réponse.
- `llm.with_structured_output(Modèle)` renvoie une instance typée et validée.
- Les `description=...` des champs sont injectées dans le prompt : c'est là que se joue la qualité de l'extraction.
- `Literal` et les contraintes (`ge`, `le`, `pattern`) verrouillent le format.

Suite

Au prochain cours : donner au LLM la possibilité d'**appeler des outils** (fonctions Python, APIs, recherche web).

Outils et appels de fonctions

Bibliothèques pour ce chapitre

```
# Bibliothèques à installer (chapitre 03)  
  
# Aucune nouvelle dépendance : '@tool' et 'bind_tools' viennent de langchain-core.  
# Tout est déjà disponible après 'uv sync'.
```

@tool et bind_tools viennent de langchain-core.

Pourquoi des outils ?

Un LLM seul a des **angles morts** bien connus :

- Pas d'arithmétique fiable au-delà de quelques chiffres.
- Pas de date du jour, pas d'heure, pas d'accès au système.
- Pas de recherche web, pas de base de données, pas d'API.
- Connaissances figées au moment de l'entraînement.

Solution

Lui exposer des **fonctions Python** qu'il peut décider d'appeler. Le LLM produit l'appel ; nous l'exécutons ; nous lui renvoyons le résultat.

Une fonction Python annotée + une docstring suffisent à devenir un outil.

- Les **annotations de type** décrivent les arguments au LLM.
- La **docstring** explique *quand* l'utiliser.
- Le nom de la fonction devient le nom de l'outil.

Bonne pratique

Une bonne docstring d'outil = une description claire *de l'usage*, pas une explication de l'implémentation.

Premier outil : voir la demande d'appel

```
1 """Premier outil : exposer une fonction Python au LLM."""
2 import os
3 from dotenv import load_dotenv
4 from langchain_core.tools import tool
5
6 load_dotenv()
7 backend = os.getenv("LLM_BACKEND", "mistral")
8
9 if backend == "mistral":
10     from langchain_mistralai import ChatMistralAI
11     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
12 else:
13     from langchain_ollama import ChatOllama
14     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
15
16
17 @tool
18 def multiplier(a: int, b: int) -> int:
19     """Multiplie deux entiers."""
20     return a * b
21
22
23 llm_avec_outil = llm.bind_tools([multiplier])
24 reponse = llm_avec_outil.invoke("Combien font 137 fois 89 ?")
25
26 print("Contenu :", reponse.content or "(vide : le modèle a appelé un outil)")
27 print("tool_calls :", reponse.tool_calls)
```

Que se passe-t-il ?

- 1 `llm.bind_tools([multiplier])` : le LLM sait désormais qu'il a accès à `multiplier(a:int, b:int)`.
- 2 Sur « Combien font 137 fois 89 ? », le modèle ne calcule pas : il **demande** d'appeler `multiplier(a=137, b=89)`.
- 3 `reponse.content` est vide ; `reponse.tool_calls` contient la demande d'appel.

À retenir

`bind_tools` ne déclenche **aucune exécution** : le LLM décide juste *quoi* appeler. C'est à nous d'exécuter et de lui rendre la main.

Cycle complet : appel, exécution, réponse

```
1 """Cycle complet : LLM demande l'outil, on l'exécute, on renvoie le résultat."""
2 import os
3 from dotenv import load_dotenv
4 from langchain_core.tools import tool
5 from langchain_core.messages import HumanMessage, ToolMessage
6
7 load_dotenv()
8 backend = os.getenv("LLM_BACKEND", "mistral")
9
10 if backend == "mistral":
11     from langchain_mistralai import ChatMistralAI
12     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
13 else:
14     from langchain_ollama import ChatOllama
15     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
16
17 @tool
18 def multiplier(a: int, b: int) -> int:
19     """Multiplie deux entiers."""
20     return a * b
21
22
23
24 outils = {"multiplier": multiplier}
25 llm_avec_outils = llm.bind_tools(list(outils.values()))
26
27 messages = [HumanMessage("Combien font 137 fois 89 ?")]
28 ai_msg = llm_avec_outils.invoke(messages)
29 messages.append(ai_msg)
30
31 for call in ai_msg.tool_calls:
32     resultat = outils[call["name"]].invoke(call["args"])
33     messages.append(ToolMessage(content=str(resultat), tool_call_id=call["id"]))
34
35 reponse_finale = llm_avec_outils.invoke(messages)
36 print(reponse_finale.content)
```

Anatomie du cycle

- 1 HumanMessage : la question initiale.
- 2 Premier invoke : le LLM renvoie un AIMessage contenant des tool_calls.
- 3 Pour chaque appel, on exécute l'outil et on ajoute un ToolMessage avec le résultat (lié par tool_call_id).
- 4 Second invoke : le LLM voit le résultat et formule la réponse en langage naturel.

Ce cycle est le cœur d'un agent. Au prochain chapitre, create_agent l'automatisera, mais le fonctionnement reste exactement celui-ci.

Plusieurs outils : le LLM choisit

```
1 """Plusieurs outils : le LLM choisit le bon selon la question."""
2 import os
3 from datetime import date
4 from dotenv import load_dotenv
5 from langchain_core.tools import tool
6 from langchain_core.messages import HumanMessage, ToolMessage
7
8 load_dotenv()
9 backend = os.getenv("LLM_BACKEND", "mistral")
10
11 if backend == "mistral":
12     from langchain_mistralai import ChatMistralAI
13     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
14 else:
15     from langchain_ollama import ChatOllama
16     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
17
18 @tool
19 def multiplier(a: int, b: int) -> int:
20     """Multiplie deux entiers."""
21     return a * b
22
23
24 @tool
25 def date_du_jour() -> str:
26     """Renvoie la date du jour au format ISO (AAAA-MM-JJ)."""
27     return date.today().isoformat()
28
29
30 outils = {"multiplier": multiplier, "date_du_jour": date_du_jour}
31 llm_avec_outils = llm.bind_tools(list(outils.values()))
32
33 for question in ["Quelle est la date du jour ?", "Calcule 42 fois 17."]:
34     messages = [HumanMessage(question)]
35     ai_msg = llm_avec_outils.invoke(messages)
36     messages.append(ai_msg)
37     for call in ai_msg.tool_calls:
38         res = outils[call["name"]].invoke(call["args"])
39         messages.append(ToolMessage(content=str(res), tool_call_id=call["id"]))
40     reponse = llm_avec_outils.invoke(messages)
41     print(f"Q: {question}")
42     print(f"R: {reponse.content}\n")
43
```

Comment le modèle choisit-il ?

- Il reçoit la liste complète des outils avec leurs signatures et docstrings.
- Il choisit en fonction du **contenu de la question** : « date du jour » déclenche `date_du_jour`, « calcule » déclenche `multiplier`.
- Il peut aussi décider de **ne rien appeler** et répondre directement.

Bonne pratique

Donner peu d'outils, bien décrits, vaut mieux qu'un catalogue mal documenté. Le LLM se trompe d'autant plus qu'il y a d'options ambiguës.

Appels parallèles

Pour une question portant sur plusieurs sujets indépendants, le LLM peut renvoyer **plusieurs tool_calls** en une seule réponse.

```
1 """Appels parallèles : le LLM peut demander plusieurs outils en une seule réponse."""
2 import os
3 from dotenv import load_dotenv
4 from langchain_core.tools import tool
5 from langchain_core.messages import HumanMessage, ToolMessage
6
7 load_dotenv()
8 backend = os.getenv("LLM_BACKEND", "mistral")
9
10 if backend == "mistral":
11     from langchain_mistralai import ChatMistralAI
12     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
13 else:
14     from langchain_ollama import ChatOllama
15     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
16
17 @tool
18 def temperature(ville: str) -> int:
19     """Renvoie la température en Celsius pour une ville (donnée fictive)."""
20     return {"Paris": 18, "Tokyo": 24, "Lima": 22}.get(ville, 20)
21
22
23
24 llm_avec_outil = llm.bind_tools([temperature])
25 messages = [HumanMessage("Quelles sont les températures à Paris, Tokyo et Lima ?")]
26 ai_msg = llm_avec_outil.invoke(messages)
27 messages.append(ai_msg)
28
29 print(f"Le modèle a fait {len(ai_msg.tool_calls)} appel(s) :")
30 for call in ai_msg.tool_calls:
31     res = temperature.invoke(call["args"])
32     print(f" - {call['name']}({call['args']}) = {res}")
33     messages.append(ToolMessage(content=str(res), tool_call_id=call["id"]))
34
35 reponse = llm_avec_outil.invoke(messages)
36 print(f"\nRéponse : {reponse.content}")
```

- Trois appels `temperature` en parallèle = une seule itération LLM au lieu de trois.
- Latence et coût divisés par autant.
- Tous les fournisseurs majeurs supportent ces appels parallèles (Mistral, OpenAI, Anthropic, Ollama avec modèles récents).

- ➊ Ajoutez un outil `additionner(a:int, b:int) -> int` et testez « Combien font $12 + 7$ puis multiplié par 3? ». Combien d'allers-retours observe-t-on?
- ➋ Modifiez `04_appels_paralleles.py` pour exécuter les outils **en parallèle** avec `concurrent.futures`.
- ➌ Faites un outil `rechercher_wikipedia(sujet:str) -> str` (utilisez la lib `wikipedia`). Demandez « Qui est Alan Turing? ».
- ➍ Que se passe-t-il si vous donnez deux outils aux noms quasi identiques (`temperature_celsius` et `temperature_kelvin`)?

Récapitulatif

- Le décorateur `@tool` transforme une fonction Python en outil exposable au LLM.
- `llm.bind_tools([...])` liste les outils disponibles ; aucune exécution.
- Le cycle **appel** → **exécution** → **réponse** se gère manuellement avec `HumanMessage`, `AIMessage`, `ToolMessage`.
- Le LLM peut demander **plusieurs** appels en une fois (parallèles).

Suite

Au prochain cours : `create_agent` automatise la boucle et nous fait passer de « LLM avec outils » à **agent** au sens propre.

Premier agent avec `create_agent`

Bibliothèques pour ce chapitre

```
# Bibliothèques à installer (chapitre 04)  
  
# Aucune nouvelle dépendance : 'create_agent' vient de langchain.agents (langchain >= 1.0).  
# Tout est déjà disponible après 'uv sync'.
```

`create_agent` vient de `langchain.agents` (`langchain >= 1.0`).

Au chapitre 3, on a écrit **à la main** la boucle :

- 1 Premier `invoke` : le LLM produit des `tool_calls`.
- 2 Pour chaque appel, on exécute l'outil et on construit un `ToolMessage`.
- 3 Second `invoke` : le LLM voit les résultats et formule la réponse.

Limite de la boucle manuelle

Si le LLM a besoin de **deux étapes successives** (ex. : la température puis une soustraction), il faut écrire la boucle *tant que* `tool_calls` n'est pas vide. C'est ce que `create_agent` fait pour nous.

Un agent, c'est une boucle automatisée

LLM → tool_calls ? → exécuter → rappeler le LLM → ...
...→ pas de tool_call → réponse finale.

- L'agent boucle **tant que** le LLM demande des outils.
- Il s'arrête quand le LLM répond directement (sans tool_calls).
- Il garde tout l'historique : chaque tour voit *tous* les messages précédents.

Premier agent

```
1 """Premier agent : create_agent automatise toute la boucle outil/exec/réponse."""
2 import os
3 from dotenv import load_dotenv
4 from langchain.agents import create_agent
5 from langchain_core.tools import tool
6
7 load_dotenv()
8 backend = os.getenv("LLM_BACKEND", "mistral")
9
10 if backend == "mistral":
11     from langchain_mistralai import ChatMistralAI
12     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
13 else:
14     from langchain_ollama import ChatOllama
15     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
16
17 @tool
18 def multiplier(a: int, b: int) -> int:
19     """Multiplie deux entiers."""
20     return a * b
21
22
23
24 agent = create_agent(model=llm, tools=[multiplier])
25
26 resultat = agent.invoke(
27     {"messages": [{"role": "user", "content": "Combien font 137 fois 89 ?"}]}
28 )
29 print(resultat["messages"][-1].content)
```

Que se passe-t-il sous le capot ?

- 1 `create_agent(model=llm, tools=[multiplier])` construit un **graphe** (LangGraph) à deux nœuds : *model* et *tools*.
- 2 `agent.invoke({"messages": [...]})` démarre la boucle : alternance `model` → `tools` → `model` → ...
- 3 Le résultat est un dict avec une clé `messages` contenant tout l'historique.
- 4 `resultat["messages"][-1]` est le dernier AIMessage : la **réponse finale**.

À retenir

Toute la boucle manuelle du chapitre 3 tient maintenant en **trois lignes**.

Donner un rôle : system_prompt

Le system_prompt est injecté en tête de chaque conversation. Idéal pour cadrer le ton, le format, ou rappeler quels outils privilégier.

```
1 """system_prompt : donner un rôle ou des consignes globales à l'agent."""
2 import os
3 from datetime import date
4 from dotenv import load_dotenv
5 from langchain.agents import create_agent
6 from langchain_core.tools import tool
7
8 load_dotenv()
9 backend = os.getenv("LLM_BACKEND", "mistral")
10
11 if backend == "mistral":
12     from langchain_mistralai import ChatMistralAI
13     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
14 else:
15     from langchain_ollama import ChatOllama
16     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
17
18
19 @tool
20 def date_du_jour() -> str:
21     """Renvoie la date du jour au format ISO (AAAA-MM-JJ)."""
22     return date.today().isoformat()
23
24
25 agent = create_agent(
26     model=llm,
27     tools=[date_du_jour],
28     system_prompt=(
29         "Tu es un assistant concis qui répond en deux phrases maximum. "
30         "Pour toute question impliquant la date, utilise l'outil date_du_jour."
31     ),
32 )
33
34 resultat = agent.invoke(
35     {"messages": [{"role": "user", "content": "On est quel jour ?"}]}
36 )
37 print(resultat["messages"][-1].content)
```

Quand utiliser `system_prompt` ?

- Persona / ton (« réponds comme un avocat », « sois concis »).
- Format imposé (« toujours en JSON », « en deux phrases maximum »).
- Politique d'usage des outils (« n'invente jamais une date sans appeler `date_du_jour` »).
- Garde-fous (« si la question sort de ton domaine, refuse poliment »).

Bonne pratique

Plus le `system_prompt` est *spécifique*, plus l'agent est *prévisible*. Mais attention : un système trop long alourdit chaque appel.

Voir la trajectoire avec agent.stream

```
1 """Voir la trajectoire : agent.stream affiche chaque message de la boucle."""
2 import os
3 from dotenv import load_dotenv
4 from langchain.agents import create_agent
5 from langchain_core.tools import tool
6 from langchain_core.messages import HumanMessage, AIMessage, ToolMessage
7
8 load_dotenv()
9 backend = os.getenv("LLM_BACKEND", "mistral")
10
11 if backend == "mistral":
12     from langchain_mistralai import ChatMistralAI
13     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
14 else:
15     from langchain_ollama import ChatOllama
16     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
17
18 @tool
19 def multiplier(a: int, b: int) -> int:
20     """Multiplie deux entiers."""
21     return a * b
22
23
24
25 agent = create_agent(model=llm, tools=[multiplier])
26
27 for etat in agent.stream(
28     {"messages": [{"role": "user", "content": "Combien font 23 fois 17 ?"}]},
29     stream_mode="values",
30 ):
31     msg = etat["messages"][-1]
32     if isinstance(msg, HumanMessage):
33         print(f"USER : {msg.content}")
34     elif isinstance(msg, AIMessage):
35         if msg.tool_calls:
36             calls = [f"{'name'}({tc['args']})" for tc in msg.tool_calls]
37             print(f"AGENT : appelle {calls}")
38         else:
39             print(f"AGENT : {msg.content}")
40     elif isinstance(msg, ToolMessage):
41         print(f"OUTIL : {msg.content}")
```

Lecture d'une trajectoire

Une trajectoire typique avec un seul outil :

- 1 HumanMessage : la question de l'utilisateur.
- 2 AIMessage avec `tool_calls` non vide : l'agent demande l'outil.
- 3 ToolMessage : le résultat brut de l'exécution.
- 4 AIMessage sans `tool_calls` : la réponse finale.

Plus la question est composée, plus la trajectoire s'allonge (étape 2 et 3 répétées).

`stream_mode="values"` renvoie l'état complet à chaque pas. Pratique pour debug ; on peut aussi capturer les `tool_calls` pour journalisation.

Multi-outils enchaînés

L'agent décide **tout seul** dans quel ordre appeler les outils, et combien de fois.

```
1 """Agent multi-outils : enchaîner plusieurs appels pour une question composée."""
2 import os
3 from datetime import date
4 from dotenv import load_dotenv
5 from langchain.agents import create_agent
6 from langchain_core.tools import tool
7
8 load_dotenv()
9 backend = os.getenv("LLM_BACKEND", "mistral")
10
11 if backend == "mistral":
12     from langchain_mistralai import ChatMistralAI
13     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
14 else:
15     from langchain_ollama import ChatOllama
16     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
17
18 @tool
19 def temperature(ville: str) -> int:
20     """Renvoie la température en Celsius pour une ville (donnée fictive)."""
21     return {"Paris": 18, "Tokyo": 24, "Lima": 22}.get(ville, 20)
22
23
24
25 @tool
26 def soustraire(a: int, b: int) -> int:
27     """Renvoie a - b."""
28     return a - b
29
30
31 @tool
32 def date_du_jour() -> str:
33     """Renvoie la date du jour au format ISO."""
34     return date.today().isoformat()
35
36
37 agent = create_agent(model=llm, tools=[temperature, soustraire, date_du_jour])
38
39 question = (
40     "Quel est l'écart de température entre Tokyo et Paris aujourd'hui, "
41     "et quelle est la date du jour ?"
42 )
43 resultat = agent.invoke({"messages": [{"role": "user", "content": question}]}
44 print(resultat["messages"][-1].content)
```

- Le LLM peut **boucler** : prévoir une limite (`recursion_limit` dans config).
- Plus il y a d'outils, plus le LLM se trompe : viser **moins mais mieux décrit**.
- Coût : chaque tour de boucle = un appel LLM. Mistral « small » reste bon marché ; pour des modèles premium, surveiller.
- Sécurité : les outils sont du code Python qui s'exécute. Ne jamais exposer un `eval` ou un accès filesystem à un agent ouvert au public.

- 1 Comparez ligne à ligne `02_executer_outil.py` (chapitre 3) et `01_premier_agent.py` : que disparaît-il ?
- 2 Modifiez `02_system_prompt.py` pour interdire au modèle de répondre sur la météo. Vérifiez en lui posant la question.
- 3 Dans `03_trajectoire.py`, ajoutez un compteur de tours et arrêtez la boucle après 5 messages, même si l'agent n'a pas terminé.
- 4 Reprenez l'exercice 3 du chapitre 3 (Wikipedia) et écrivez l'équivalent avec `create_agent`. Comparez la concision du code.

Récapitulatif

- `create_agent(model, tools)` construit un agent qui boucle automatiquement LLM ↔ tools.
- `system_prompt` cadre le comportement global de l'agent.
- `agent.stream(..., stream_mode="values")` expose la trajectoire pour debug ou journalisation.
- L'agent enchaîne plusieurs outils sans qu'on ait à scripter l'ordre.

Suite

Au prochain cours : doter l'agent d'une **mémoire** entre conversations, et **streamer** ses tokens en temps réel.

Mémoire et streaming

```
# Bibliothèques à installer (chapitre 05)
```

```
# 'langgraph' est installé via langchain (transitive dependency).
```

```
# Si besoin explicite :
```

```
uv add langgraph
```

Pourquoi mémoire et streaming ?

Jusqu'ici, chaque agent .invoke est **indépendant** : aucune mémoire d'un appel à l'autre.

- Pas de « Comment je m'appelle ? » après un « Je m'appelle Marie » dans un appel précédent.
- Pas de chatbot multi-tour, pas d'assistant qui apprend de l'utilisateur dans la session.

Et la sortie arrive d'un **seul bloc** après plusieurs secondes :

- L'utilisateur attend dans le vide.
- Pour des longues réponses (RAG, rapports...), c'est inacceptable en UX.

Solution : un **checkpointer** pour la mémoire, un **stream** pour les tokens.

Le checkpointer : sauvegarder l'état entre appels

- Un checkpointer est un magasin clé/valeur où l'agent sérialise tout son état après chaque pas.
- La clé est le `thread_id` : un identifiant de conversation.
- Au tour suivant, l'agent recharge l'état du thread, ajoute le nouveau message et continue.

Trois implémentations courantes :

- `InMemorySaver` (langgraph) : en RAM, disparaît à l'arrêt du process. Idéal pour développer.
- `PostgresSaver`, `SqliteSaver` : persistance disque, durable.
- `RedisSaver` : partage entre instances, latence faible.

Mémoire entre tours

```
1 """Mémoire entre tours : checkpointer + thread_id."""
2 import os
3 from dotenv import load_dotenv
4 from langchain.agents import create_agent
5 from langgraph.checkpoint.memory import InMemorySaver
6
7 load_dotenv()
8 backend = os.getenv("LLM_BACKEND", "mistral")
9
10 if backend == "mistral":
11     from langchain_mistralai import ChatMistralAI
12     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
13 else:
14     from langchain_ollama import ChatOllama
15     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
16
17 agent = create_agent(model=llm, tools=[], checkpointer=InMemorySaver())
18 config = {"configurable": {"thread_id": "session-1"}}
19
20
21 r1 = agent.invoke({"messages": [{"role": "user", "content": "Bonjour, je m'appelle Marie."}]}, config)
22 print("Tour 1 :", r1["messages"][-1].content)
23
24 r2 = agent.invoke({"messages": [{"role": "user", "content": "Comment je m'appelle ?"}]}, config)
25 print("Tour 2 :", r2["messages"][-1].content)
```

- 1 `create_agent(..., checkpointer=InMemorySaver())` : l'agent saura sauvegarder son état.
- 2 Le `config={"configurable": {"thread_id": "session-1"}}` identifie la conversation.
- 3 Au premier `invoke`, l'historique contient [Marie]. Au second, l'agent voit *tout* l'historique et peut répondre « Vous vous appelez Marie ».

À retenir

Sans `thread_id`, le `checkerpoint` ne peut rien sauvegarder. Avec un `thread_id` fixe, on construit une **conversation continue**.

Threads séparés = conversations isolées

```
1 """Deux threads = deux conversations isolées."""
2 import os
3 from dotenv import load_dotenv
4 from langchain.agents import create_agent
5 from langgraph.checkpoint.memory import InMemorySaver
6
7 load_dotenv()
8 backend = os.getenv("LLM_BACKEND", "mistral")
9
10 if backend == "mistral":
11     from langchain_mistralai import ChatMistralAI
12     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
13 else:
14     from langchain_ollama import ChatOllama
15     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
16
17 agent = create_agent(model=llm, tools=[], checkpointer=InMemorySaver())
18
19 agent.invoke(
20     {"messages": [{"role": "user", "content": "Je m'appelle Bob."}],
21      "configurable": {"thread_id": "alice"}},
22 )
23 agent.invoke(
24     {"messages": [{"role": "user", "content": "Je m'appelle Alice."}],
25      "configurable": {"thread_id": "alice"}},
26 )
27
28 # Sur le thread "bob", l'agent ne sait rien d'Alice ni de Bob.
29 r = agent.invoke(
30     {"messages": [{"role": "user", "content": "Tu connais mon prénom ?"}],
31      "configurable": {"thread_id": "bob"}},
32 )
33
34 print("Thread bob :", r["messages"][-1].content)
```

Pourquoi des threads ?

- Un même agent dessert **plusieurs utilisateurs** : un `thread_id` par session web ou par client API.
- Aucune fuite d'information d'un thread à l'autre : c'est l'isolation par défaut.
- Un même utilisateur peut avoir plusieurs conversations parallèles (ex : « conversation projet A », « conversation projet B »).

Bonne pratique

Choisir un `thread_id` stable et secret côté serveur (UUID, hash de session...). Ne jamais le passer en clair côté client.

Streaming token par token

Au lieu d'attendre la réponse complète, on consomme les tokens **au fur et à mesure**.

```
1 """Streaming token par token : stream_mode="messages"."""
2 import os
3 from dotenv import load_dotenv
4 from langchain.agents import create_agent
5 from langchain_core.messages import AIMessageChunk
6
7 load_dotenv()
8 backend = os.getenv("LLM_BACKEND", "mistral")
9
10 if backend == "mistral":
11     from langchain_mistralai import ChatMistralAI
12     llm = ChatMistralAI(model="mistral-small-latest", temperature=0)
13 else:
14     from langchain_ollama import ChatOllama
15     llm = ChatOllama(model=os.getenv("OLLAMA_MODEL", "qwen2.5:3b"), temperature=0)
16
17 agent = create_agent(model=llm, tools=[])
18
19 for chunk in agent.stream(
20     {"messages": [{"role": "user", "content": "Donne-moi trois conseils pour apprendre Python."}]},
21     stream_mode="messages",
22     version="v2",
23 ):
24     if chunk["type"] == "messages":
25         token, _ = chunk["data"]
26         if isinstance(token, AIMessageChunk) and token.text:
27             print(token.text, end="", flush=True)
28
29 print()
```

Modes de streaming

`agent.stream(..., stream_mode=X, version="v2")` accepte plusieurs modes :

- "messages" : tokens individuels (`AIMessageChunk`). Idéal pour afficher en direct.
- "values" : état complet à chaque pas (*vu au chapitre 4*). Idéal pour debug.
- "updates" : seulement les changements (sortie des nœuds). Idéal pour journalisation.
- Liste de modes : `["messages", "updates"]` pour combiner.

Avec `version="v2"`, chaque chunk a une clé `type` et une clé `data` : on peut router selon le type. C'est le format moderne.

Refactoring : centraliser le choix du backend

Tous nos snippets dupliquent le bloc `if backend == "mistral"...`. Le helper `code/_lib/llm.py` centralise cette logique.

```
1  """Refactoring : extraire la sélection de backend dans '_lib/llm.py'.
```

```
2
```

```
3  Le helper 'get_llm()' renvoie ChatMistralAI ou ChatOllama selon LLM_BACKEND.
```

```
4  On ajoute 'code/' à sys.path pour que 'from _lib.llm import get_llm' fonctionne
```

```
5  quel que soit le dossier de lancement.
```

```
6  """
```

```
7  import sys
```

```
8  from pathlib import Path
```

```
9  sys.path.insert(0, str(Path(__file__).parent.parent))
```

```
10
```

```
11 from dotenv import load_dotenv
```

```
12 from langchain.agents import create_agent
```

```
13 from langchain_core.tools import tool
```

```
14 from _lib.llm import get_llm
```

```
15
```

```
16 load_dotenv()
```

```
17 llm = get_llm(temperature=0)
```

```
18
```

```
19
```

```
20 @tool
```

```
21 def multiplier(a: int, b: int) -> int:
```

```
22     """Multiplie deux entiers."""
```

```
23     return a * b
```

```
24
```

```
25
```

```
26 agent = create_agent(model=llm, tools=[multiplier])
```

```
27 resultat = agent.invoke(
```

```
28     {"messages": [{"role": "user", "content": "Combien font 12 fois 11 ?"}]}
```

```
29 )
```

```
30 print(resultat["messages"][-1].content)
```

Pourquoi cette refactorisation maintenant ?

- Au chapitre 1, l'instanciation directe avait une vertu pédagogique (« voir le ChatModel choisi »).
- Avec 4-5 snippets par chapitre, la duplication devient bruyante et masque l'intention de chaque snippet.
- Centraliser permet de changer de modèle ou ajouter un backend sans toucher 20 fichiers.

Astuce Python

`sys.path.insert(0, ...)` est nécessaire car les snippets sont exécutés depuis leur dossier (pas depuis la racine). Pour un *vrai* projet, on transformerait `_lib` en package installable (`pyproject.toml`) ou on utiliserait `src/-layout`.

- 1 Modifiez `01_memoire_simple.py` pour ajouter un troisième tour qui demande le carré de l'âge mentionné au tour 1.
- 2 Dans `02_threads_separés.py`, ajoutez un troisième thread qui parle d'un sujet sans lien (recette de cuisine). Vérifiez l'isolation.
- 3 Combinez `stream_mode=["messages", "updates"]` dans `03_streaming_tokens.py` pour afficher **à la fois** les tokens et les appels d'outils.
- 4 Refactorisez les snippets du chapitre 4 pour utiliser `get_llm()`. Combien de lignes économisées ?

Récapitulatif

- `checkpointer=InMemorySaver()` + `thread_id` dans le config = mémoire entre appels.
- Plusieurs `thread_id` = conversations isolées en parallèle.
- `stream_mode="messages"` avec `version="v2"` = tokens en direct.
- Centraliser la sélection de backend dans `_lib/llm.py` simplifie l'évolution du cours.

Suite

Au prochain cours : **LangGraph** pour les boucles non triviales (validation humaine, branchements conditionnels, sous-graphes).

LangGraph pour les boucles complexes

```
# Bibliothèques à installer (chapitre 06)  
  
# 'langgraph' est déjà installé (transitive via langchain).  
# Si besoin explicite :  
uv add langgraph
```

Pourquoi sortir de `create_agent` ?

`create_agent` est parfait pour la boucle **LLM ↔ outils**. Mais beaucoup de cas réels demandent davantage :

- Branchement **conditionnel** : si le message est du spam, ne pas appeler le LLM cher.
- **Validation humaine** avant l'exécution d'un outil sensible.
- Pipelines multi-étapes : classifier → raffiner → valider → répondre.
- Boucles avec **condition d'arrêt métier** (jusqu'à ce qu'un score atteigne X).

Bonne nouvelle

`create_agent` **est** un `StateGraph` en interne. On apprend ici à le construire à la main.

Anatomie d'un StateGraph

Trois ingrédients seulement :

- **State** : un TypedDict (ou Pydantic BaseModel) décrivant ce qui circule.
- **Nodes** : des fonctions `state → partial state`. Pures, déterministes, ou appelant un LLM.
- **Edges** : qui pointe vers qui. Statiques ou conditionnelles.

Deux nœuds spéciaux : START (point d'entrée) et END (sortie).

Ce qui change

Avec `create_agent`, on *configure* un graphe préfait. Avec `StateGraph`, on **décrit** le graphe nous-mêmes.

Premier graphe : pipeline linéaire

```
1  """Premier StateGraph : un état, deux nœuds, deux arêtes."""
2  from typing import TypedDict
3  from langgraph.graph import StateGraph, START, END
4
5
6  class Etat(TypedDict):
7      sujet: str
8      blague: str | None
9
10
11 def raffiner_sujet(etat: Etat) -> Etat:
12     return {"sujet": etat["sujet"] + " et les chats"}
13
14
15 def generer_blague(etat: Etat) -> Etat:
16     return {"blague": f"Voici une blague sur {etat['sujet']}..."}
17
18
19 graphe = (
20     StateGraph(Etat)
21     .add_node(raffiner_sujet)
22     .add_node(generer_blague)
23     .add_edge(START, "raffiner_sujet")
24     .add_edge("raffiner_sujet", "generer_blague")
25     .add_edge("generer_blague", END)
26     .compile()
27 )
28
29 resultat = graphe.invoke({"sujet": "les pompiers", "blague": None})
30 print(resultat)
```

- 1 Etat déclare les champs persistants (sujet, blague).
- 2 Chaque nœud retourne un **dict partiel** : seules les clés présentes sont mises à jour.
- 3 `add_edge(START, "raffiner_sujet")` démarre par ce nœud ;
`add_edge("generer_blague", END)` termine.
- 4 `compile()` produit un **Pregel** qui s'invoque comme un agent.

Branchement conditionnel

Une fonction router retourne le nom du prochain nœud (ou END).

```
1 """Branchement conditionnel : route le flux selon l'état."""
2 from typing import Literal, TypedDict
3 from langgraph.graph import StateGraph, START, END
4
5
6 class Etat(TypedDict):
7     message: str
8     longueur: int | None
9     reponse: str | None
10
11
12 def mesurer(etat: Etat) -> Etat:
13     return {"longueur": len(etat["message"])}
14
15
16 def reponse_courte(etat: Etat) -> Etat:
17     return {"reponse": "Message court reçu."}
18
19
20 def reponse_longue(etat: Etat) -> Etat:
21     return {"reponse": "Message long, je résume mentalement."}
22
23
24 def router(etat: Etat) -> Literal["reponse_courte", "reponse_longue"]:
25     return "reponse_courte" if etat["longueur"] < 30 else "reponse_longue"
26
27
28 graphe = (
29     StateGraph(Etat)
30     .add_node(mesurer)
31     .add_node(reponse_courte)
32     .add_node(reponse_longue)
33     .add_edge(START, "mesurer")
34     .add_conditional_edges("mesurer", router)
35     .add_edge("reponse_courte", END)
36     .add_edge("reponse_longue", END)
37     .compile()
38 )
39
40 for msg in ["Salut", "Pourrais-tu me détailler en quelques lignes le principe de la photosynthèse ?"]:
41     etat = graphe.invoke({"message": msg, "longueur": None, "reponse": None})
42     print(f"({etat['longueur']} car) {etat['reponse']}")
```

- `add_conditional_edges("mesurer", router)` : après mesurer, c'est router qui décide.
- Le type de retour de `router` (`Literal["reponse_courte", "reponse_longue"]`) sert de documentation et de validation.
- Toutes les destinations possibles doivent être des nœuds existants ou `END`.

Bonne pratique

Garder le routeur **idempotent et rapide** : ne pas y appeler de LLM, juste lire l'état.

Boucle avec condition d'arrêt

Le Annotated[... , operator.add] indique à LangGraph que les retours doivent être **accumulés** plutôt qu'écrasés.

```
1  """Boucle avec arrêt conditionnel : itérer jusqu'à atteindre une condition."""
2  import operator
3  from typing import Annotated, Literal, TypedDict
4  from langgraph.graph import StateGraph, START, END
5
6
7  class Etat(TypedDict):
8      # Le reducer 'operator.add' accumule au lieu d'écraser.
9      valeurs: Annotated[list[int], operator.add]
10
11
12  def ajouter(etat: Etat) -> Etat:
13      n = len(etat["valeurs"])
14      return {"valeurs": [n + 1]}
15
16
17  def continuer(etat: Etat) -> Literal["ajouter", END]:
18      return "ajouter" if sum(etat["valeurs"]) < 20 else END
19
20
21  graphe = (
22      StateGraph(Etat)
23      .add_node(ajouter)
24      .add_edge(START, "ajouter")
25      .add_conditional_edges("ajouter", continuer)
26      .compile()
27  )
28
29  resultat = graphe.invoke({"valeurs": []})
30  print("Valeurs :", resultat["valeurs"])
31  print("Somme   :", sum(resultat["valeurs"]))
```

Sans reducer, chaque nœud **remplace** la valeur précédente. Avec un reducer, il **fusionne**.

- `operator.add` sur une liste : concatène.
- `add_messages` (`langgraph`) : ajoute des messages en évitant les doublons (utilisé par `create_agent`).
- Un reducer custom est juste une fonction (`old, new`) \rightarrow `merged`.

Une boucle dans `LangGraph` est une simple arête conditionnelle qui peut renvoyer vers un nœud déjà visité. Pas de `while` explicite.

Hybride : classifieur puis expert (LLM)

Un graphe peut mélanger nœuds Python et nœuds appelant un LLM. Ici un classifieur LLM route vers le bon expert LLM.

```
1 """Hybride LLM : un classifieur route vers le bon expert."""
2 import sys
3 from pathlib import Path
4 sys.path.insert(0, str(Path(__file__).parent.parent))
5
6 from typing import Literal, TypedDict
7 from dotenv import load_dotenv
8 from langgraph.graph import StateGraph, START, END
9 from lib.llm import get_llm
10
11 load_dotenv()
12 llm = get_llm(temperature=0)
13
14
15 class Etat(TypedDict):
16     question: str
17     domaine: str | None
18     reponse: str | None
19
20
21 def classifieur(etat: Etat) -> Etat:
22     prompt = (
23         f"La question suivante porte-t-elle sur les mathématiques ou autre chose ?"
24         f"Réponds par un seul mot parmi {{math, autre}}.\n\nQuestion : {etat['question']}"
25     )
26     domaine = llm.invoke(prompt).contenu.strip().lower()
27     return {"domaine": "math" if "math" in domaine else "autre"}
28
29
30 def expert_math(etat: Etat) -> Etat:
31     r = llm.invoke(f"Tu es prof de maths. Réponds en une phrase.\n\n{etat['question']}")
32     return {"reponse": r.contenu}
33
34
35 def expert_general(etat: Etat) -> Etat:
36     r = llm.invoke(f"Réponds brièvement.\n\n{etat['question']}")
37     return {"reponse": r.contenu}
38
39
40 def router(etat: Etat) -> Literal["expert_math", "expert_general"]:
41     return "expert_math" if etat["domaine"] == "math" else "expert_general"
42
43
44 graphe = (
45     StateGraph(Etat)
46     .add_node(classifieur)
47     .add_node(expert_math)
48     .add_node(expert_general)
49     .add_edge(START, "classifieur")
50     .add_conditional_edges("classifieur", router)
51     .add_edge("expert_math", END)
52     .add_edge("expert_general", END)
53     .compile()
54 )
55
56 for q in ["Combien fait 12 fois 13 ?", "Quel est le plus grand fleuve d'Afrique ?"]:
57     etat = graphe.invoke({"question": q, "domaine": None, "reponse": None})
58     print(f"[{etat['domaine']}]: Q: {q}")
59     print(f"      R: {etat['reponse']}\n")
```

Pourquoi diviser en nœuds ?

- Chaque nœud a un **prompt court et spécifique** : meilleur résultat qu'un seul méga-prompt.
- Le `system_prompt` de chaque expert peut être adapté à son domaine.
- On peut surveiller, logger, mettre en cache au niveau d'un seul nœud.
- On peut **remplacer** un nœud (par ex. `expert_math` → outil calculatrice) sans toucher au reste.

create_agent vs StateGraph

- create_agent : 90 % des cas, **une ligne** suffit. Bon défaut.
- StateGraph : quand on a besoin de **contrôle fin** (validation humaine, branchements métier, sous-graphes, multi-agents).
- Sous le capot, create_agent est un StateGraph avec deux nœuds (model, tools) et une boucle conditionnelle.

Astuce

On peut visualiser un graphe avec `graphe.get_graph().draw_mermaid_png()`. Pratique en debug.

- 1 Modifiez `02_branchement.py` pour ajouter un troisième cas « moyen » (entre 30 et 100 caractères).
- 2 Dans `03_boucle.py`, faites afficher un message à chaque tour de boucle (ajoutez un nœud tracer).
- 3 Étendez `04_classifier_expert.py` avec un troisième expert (« histoire »). Ajoutez les arêtes nécessaires.
- 4 Reconstituez à la main `create_agent` avec un `StateGraph` : un nœud `model`, un nœud `tools`, une boucle conditionnelle. Bonus : faites-le boucler tant que `tool_calls` n'est pas vide.

Récapitulatif

- `StateGraph(State) + add_node + add_edge + compile` = un graphe de calcul typé.
- `add_conditional_edges` permet le routage métier ; le routeur est une fonction pure rapide.
- Les **reducers** (`Annotated[... , op]`) gouvernent la fusion des états dans les boucles.
- Nœuds Python et nœuds LLM cohabitent sans friction.

Suite

Au prochain et dernier cours : un **RAG minimal** pour donner à l'agent un accès à des documents externes (PDF, base vectorielle).

RAG minimal

```
# Bibliothèques à installer (chapitre 07)

uv add langchain-text-splitters numpy

# Si LLM_BACKEND=ollama, il faut aussi un modèle d'embeddings :
ollama pull nomic-embed-text
```

Pourquoi RAG ?

Un LLM a une **connaissance figée** (date d'entraînement) et ne sait *rien* de :

- Vos documents internes (notes, mails, manuels, contrats).
- Les évolutions postérieures à son entraînement.
- Les bases métier propres à votre organisation.

Idée

Retrieval-**A**ugmented **G**eneration : avant de répondre, on *retrouve* les passages pertinents dans une base, et on les *fournit en contexte* au LLM.

Pipeline RAG en cinq étapes

- ➊ **Loader** : charger les documents (texte, PDF, web, CSV...).
- ➋ **Splitter** : découper en *chunks* de taille gérable.
- ➌ **Embeddings** : convertir chaque chunk en vecteur dense.
- ➍ **VectorStore** : stocker les vecteurs (en mémoire, Postgres+pgvector, Qdrant...).
- ➎ **Retriever** : pour une question, retourner les k chunks les plus proches en cosinus.

Le LLM n'intervient qu'à la **toute fin**, en lisant question + chunks retrouvés.

Les embeddings, en deux mots

Un **embedding** est un vecteur dense (typiquement 768-1536 dimensions) tel que des textes *sémantiquement proches* ont des vecteurs proches.

- Mistral : `mistral-embed` (1024 dim, cloud).
- Ollama : `nomic-embed-text` (768 dim, local).
- OpenAI : `text-embedding-3-small` (1536 dim, cloud).

Bonne pratique

Toujours *indexer* et *interroger* avec le même modèle d'embeddings. Sinon les distances n'ont aucun sens.

Étape 1 : indexer un document

```
1 """Étape 1 : découper un document, calculer les embeddings, indexer en mémoire."""
2 import sys
3 from pathlib import Path
4 sys.path.insert(0, str(Path(__file__).parent.parent))
5
6 from dotenv import load_dotenv
7 from langchain_core.documents import Document
8 from langchain_core.vectorstores import InMemoryVectorStore
9 from langchain_text_splitters import RecursiveCharacterTextSplitter
10 from _lib.embeddings import get_embeddings
11
12 load_dotenv()
13 texte = (Path(__file__).parent / "data" / "iut.txt").read_text()
14
15 splitter = RecursiveCharacterTextSplitter(chunk_size=200, chunk_overlap=40)
16 chunks = splitter.split_text(texte)
17 documents = [Document(page_content=c) for c in chunks]
18
19 vectorstore = InMemoryVectorStore.from_documents(documents=documents, embedding=get_embeddings())
20 print(f"{len(documents)} chunks indexés.")
21 print("Premier chunk :", documents[0].page_content[:80], "...")
```

Choisir `chunk_size` et `chunk_overlap`

- `chunk_size` *trop petit* : on perd le contexte (« cette mesure », « cet auteur » sans antécédent).
- `chunk_size` *trop grand* : la recherche devient bruitée, le LLM doit lire plus.
- `chunk_overlap` : recouvrement entre chunks pour éviter de couper en plein milieu d'une idée.

Valeurs courantes : `chunk_size=500-1000`, `chunk_overlap=50-200`, à ajuster selon le type de texte.

Étape 2 : rechercher par similarité

```
1 """Étape 2 : recherche par similarité sémantique."""
2 import sys
3 from pathlib import Path
4 sys.path.insert(0, str(Path(__file__).parent.parent))
5
6 from dotenv import load_dotenv
7 from langchain_core.documents import Document
8 from langchain_core.vectorstores import InMemoryVectorStore
9 from langchain_text_splitters import RecursiveCharacterTextSplitter
10 from _lib.embeddings import get_embeddings
11
12 load_dotenv()
13 texte = (Path(__file__).parent / "data" / "iut.txt").read_text()
14 splitter = RecursiveCharacterTextSplitter(chunk_size=200, chunk_overlap=40)
15 documents = [Document(page_content=c) for c in splitter.split_text(texte)]
16 vectorstore = InMemoryVectorStore.from_documents(documents=documents, embedding=get_embeddings())
17
18 for question in ["Combien d'étudiants par promotion ?", "Sur quoi travaille Christophe Guyeux ?"]:
19     print(f"\nQ : {question}")
20     for doc in vectorstore.similarity_search(question, k=2):
21         print(f" - {doc.page_content[:100]}...")
```

La similarité sémantique

- `similarity_search(question, k=2)` encode la question, calcule la similarité cosinus avec tous les chunks indexés, retourne les k premiers.
- C'est **coûteux** en RAM mais simple à mettre en œuvre. Pour un cours, c'est parfait. Pour 10 millions de chunks, on passe à pgvector ou Qdrant.

Les chunks retrouvés ne sont pas toujours pertinents : on appelle ça le *ranker* et on peut le perfectionner avec un re-ranker (ex. : Cohere Rerank).

Étape 3 : RAG « naïf » en chaîne LCEL

```
1 """Étape 3 : RAG « naïf » en chaîne LCEL : retrieve → prompt → llm."""
2 import sys
3 from pathlib import Path
4 sys.path.insert(0, str(Path(__file__).parent.parent))
5
6 from dotenv import load_dotenv
7 from langchain_core.documents import Document
8 from langchain_core.prompts import ChatPromptTemplate
9 from langchain_core.output_parsers import StrOutputParser
10 from langchain_core.runnables import RunnablePassthrough
11 from langchain_core.vectorstores import InMemoryVectorStore
12 from langchain_text_splitters import RecursiveCharacterTextSplitter
13 from _lib.embeddings import get_embeddings
14 from _lib.llm import get_llm
15
16 load_dotenv()
17 texte = (Path(__file__).parent / "data" / "iut.txt").read_text()
18 splitter = RecursiveCharacterTextSplitter(chunk_size=200, chunk_overlap=40)
19 documents = [Document(page_content=c) for c in splitter.split_text(texte)]
20 retriever = InMemoryVectorStore.from_documents(documents, get_embeddings()).as_retriever(k=2)
21
22 prompt = ChatPromptTemplate.from_template(
23     "Réponds à la question en t'appuyant uniquement sur le contexte ci-dessous. "
24     "Si l'info manque, dis-le.\n\nContexte : \n{contexte}\n\nQuestion : {question}"
25 )
26
27
28 def format_docs(docs):
29     return "\n---\n".join(d.page_content for d in docs)
30
31
32 chain = (
33     {"contexte": retriever | format_docs, "question": RunnablePassthrough()}
34     | prompt
35     | get_llm(temperature=0)
36     | StrOutputParser()
37 )
38
39 print(chain.invoke("Combien d'étudiants par promotion en informatique à l'IUT NFC ?"))
```

```
{contexte: retriever, question: passthrough} → prompt → llm →  
StrOutputParser
```

- `RunnablePassthrough()` laisse passer la question telle quelle.
- Le `retriever` est appelé sur la question puis formaté en texte.
- Le LLM voit *contexte + question* et formule la réponse.

Limite

Ce RAG appelle **toujours** le retriever, même quand ce n'est pas utile (« Combien font 2 + 2? »). On améliore avec un agent.

Étape 4 : RAG agentique

L'agent décide **de lui-même** d'appeler ou non l'outil de recherche.

```
1 """Étape 4 : RAG agentique. L'agent décide quand appeler l'outil de recherche."""
2 import sys
3 from pathlib import Path
4 sys.path.insert(0, str(Path(__file__).parent.parent))
5
6 from dotenv import load_dotenv
7 from langchain.agents import create_agent
8 from langchain_core.documents import Document
9 from langchain_core.tools import tool
10 from langchain_core.vectorstores import InMemoryVectorStore
11 from langchain_text_splitters import RecursiveCharacterTextSplitter
12 from _lib.embeddings import get_embeddings
13 from _lib.llm import get_llm
14
15 load_dotenv()
16 texte = (Path(__file__).parent / "data" / "iut.txt").read_text()
17 splitter = RecursiveCharacterTextSplitter(chunk_size=200, chunk_overlap=40)
18 documents = [Document(page_content=c) for c in splitter.split_text(texte)]
19 vs = InMemoryVectorStore.from_documents(documents, get_embeddings())
20
21
22 @tool
23 def rechercher_iut(question: str) -> str:
24     """Recherche dans la documentation officielle de l'IUT NFC."""
25     return "\n--\n".join(d.page_content for d in vs.similarity_search(question, k=2))
26
27
28 agent = create_agent(model=get_llm(temperature=0), tools=[rechercher_iut])
29
30 for question in [
31     "Que fait Christophe Guyeux dans la recherche ?",
32     "Combien font 12 fois 11 ?", # pas de RAG nécessaire
33 ]:
34     r = agent.invoke({"messages": [{"role": "user", "content": question}]}
35     print(f"Q : {question}")
36     print(f"R : {r['messages'][-1].content}\n")
```

RAG naïf vs RAG agentique

- **Naïf** : 1 retrieval forcé + 1 appel LLM. Latence et coût stables, mais retrieve même quand inutile.
- **Agentique** : l'agent peut décider *ne pas* retrieve, ou retrieve *plusieurs fois* avec des reformulations. Plus souple, plus coûteux dans les cas complexes.

Tendance 2025-2026

Le RAG agentique gagne du terrain : les agents savent désormais *quand* chercher, *quoi* chercher, et *comment* reformuler après un échec.

- **Loaders riches** : PyPDFLoader, WebBaseLoader, NotionLoader, GitHubIssuesLoader.
- **VectorStores persistants** : Chroma, Qdrant, pgvector, Weaviate.
- **Re-ranking** avec un modèle dédié pour réorganiser le top- k .
- **HyDE** : générer une réponse hypothétique et faire le retrieval dessus.
- **Multi-query** : reformuler la question en plusieurs variantes pour élargir le retrieval.

- ➊ Remplacez `iut.txt` par un PDF (utilisez `PyPDFLoader`, `uv add langchain-community pypdf`). Comparez les résultats.
- ➋ Faites varier `chunk_size` entre 100 et 1000 dans `01_indexer.py`. Pour quelles questions le RAG fonctionne-t-il mieux ?
- ➌ Dans `04_rag_agentique.py`, ajoutez un second outil `multiplier` et posez une question composée « Combien d'étudiants \times 12 mois ? ». Observez la trajectoire.
- ➍ Persistez l'index sur disque avec Chroma ou FAISS pour ne pas réindexer à chaque démarrage.

Récapitulatif

- RAG = **retrieval** de chunks + **augmentation** du prompt + **generation** par le LLM.
- Pipeline : Loader → Splitter → Embeddings → VectorStore → Retriever.
- RAG naïf en LCEL : 5 lignes. RAG agentique avec `create_agent` : un outil de recherche + le LLM décide.
- Toujours indexer et interroger avec le **même modèle d'embeddings**.

Fin du cours

Vous avez maintenant les briques pour construire un **assistant documentaire** de bout en bout. À vous de jouer.

Merci.

Questions?