

# **CURSO COMPLETO DE LANGCHAIN, LANGGRAPH Y AGENTES DE IA CON PYTHON**

**SANTIAGO HERNÁNDEZ RAMOS**



# Tema 1 – Introducción a LangChain: Tu primera aplicación con IA y LLMs

## Guía de referencia

Usa esta guía como recordatorio rápido de componentes, clases y parámetros clave. Está pensada para cualquier proyecto, no solo para los ejemplos del curso.

### 1 · Configuración del entorno

Componente	Comando/Acción	Comentario
Entorno virtual	<code>python -m venv venv</code>	Crea un entorno aislado para las dependencias del proyecto
Activar entorno	<code>venv\Scripts\activate</code> (Windows)	Verifica que aparezca (venv) en tu terminal
Instalar LangChain	<code>pip install langchain langchain-openai langchain-google-genai</code>	Instala el framework base y conectores para proveedores
Estructura de carpetas	Organizar por temas/módulos	Facilita la navegación y mantenimiento del código

### 2 · Proveedores y modelos (clases y opciones)

Proveedor	Clase LangChain	Parámetros clave (ejemplos)	Notas
OpenAI	<code>from langchain_openai import ChatOpenAI</code>	<code>model="gpt-4-mini" , temperature=0.0-1.0</code>	<code>temperature</code> ↑ → más creatividad/aleatoriedad; ↓ → más determinismo/consistencia.
Google GenAI (Gemini)	<code>from langchain_google_genai import ChatGoogleGenerativeAI</code>	<code>model="gemini-2.5-flash" , temperature=...</code>	Interfaz equivalente: cambiar proveedor requiere mínimos cambios.

Patrón común de inicialización

```
llm = ChatOpenAI(model="gpt-4-mini", temperature=0.7)
# o
llm = ChatGoogleGenerativeAI(model="gemini-2.5-flash", temperature=0.7)
```

### 3 · Gestión de credenciales (variables de entorno)

Sistema	Variable típica	Ejemplo de definición	Uso en código
Windows (CMD)	<code>OPENAI_API_KEY</code>	<code>setx OPENAI_API_KEY "tu_clave"</code>	LangChain las detecta automáticamente al crear <code>ChatOpenAI</code> .
Windows (CMD)	<code>GOOGLE_API_KEY</code>	<code>setx GOOGLE_API_KEY "tu_clave"</code>	Reabre VS Code/terminal para recargar el entorno.
macOS/Linux (bash/zsh)	<code>OPENAI_API_KEY</code> / <code>GOOGLE_API_KEY</code>	<code>export OPENAI_API_KEY="tu_clave"</code>	Añádalo a tu <code>~/.bashrc</code> / <code>~/.zshrc</code> para persistir.

Buenas prácticas

- Nunca hardcodees la API key en el script.
- Cierra y reabre la terminal tras definir la variable para que el proceso herede el valor.

4 · Invocación básica del LLM (Hello World)

Paso	Qué haces	Snippet
Definir prompt	Texto de entrada (puede venir de usuario o del programa)	<code>pregunta = "¿En qué año llegó el ser humano a la Luna por primera vez?"</code>
Invocar	Usas la interfaz unificada de LangChain	<code>respuesta = llm.invoke(pregunta)</code>
Leer resultado	Obtienes el contenido generado	<code>print(respuesta.content)</code>

Notas

- `.invoke()` es la interfaz estándar (sin importar el proveedor).
- El objeto devuelto suele incluir `content` con el texto generado.

5 · Prompt Templates (plantillas reutilizables y con variables)

Componente	Función	Sintaxis
PromptTemplate	Define prompts reutilizables con variables	<code>from langchain.prompts import PromptTemplate</code>
input_variables	Lista de variables dinámicas	<code>input_variables=["nombre", "tema"]</code>
template	Texto con placeholders	<code>"Hola {nombre}, hablemos de {tema}"</code>
Variables	Se sustituyen en tiempo de ejecución	Usar <code>{variable}</code> dentro del template

Patrón de uso

```
from langchain.prompts import PromptTemplate

plantilla = PromptTemplate(
    input_variables=["nombre"],
    template="Saluda al usuario.\nNombre del usuario: {nombre}\nAsistente:"
)

# (Opcional) Render manual:
prompt_renderizado = plantilla.format(nombre="Carlos")
```

Consejos

- Usa `\n` para estructurar instrucciones (roles, pasos, secciones).
- Declara todas las variables dinámicas en `input_variables`.
- Mantén el prompt claro y específico (instrucciones → contexto → salida esperada).

6 · Cadenas (Chains) con LangChain

Objetivo: Encadenar pasos como “plantilla → LLM → salida”, con una API coherente.

Enfoque clásico (deprecado pero aún visto en código existente):

```
from langchain.chains import LLMChain

chain = LLMChain(llm=llm, prompt=plantilla)
resultado = chain.run(nombre="Carlos") # pasa las variables del prompt
print(resultado) # Texto generado por el LLM
```

Punto	Clave
LLMChain	API clásica; verás avisos de <i>deprecated</i> . Aún útil para leer/entender proyectos existentes.
Invocación	<code>.run(**kwargs)</code> pasa las variables para sustituir en la plantilla.
Evolución	La tendencia actual es usar la sintaxis basada en <i>runnables/LCEL</i> (se cubrirá en secciones posteriores).

7 · Parámetros y prácticas recomendadas

- Temperatura:
  - `0.0-0.2` → Respuestas precisas/estables (QA factual, extracción).
  - `0.3-0.7` → Equilibrio general (asistentes, reformulación).
  - `0.8-1.0` → Creatividad alta (brainstorming, copy).
- Trazabilidad: imprime/guarda `prompt` y parámetros (útil para depuración).
- Versiona prompts: tratar plantillas como artefactos (pueden romper resultados).
- Estructura por temas: carpetas por módulo/lección facilitan localizar ejemplos.
- Proveedor intercambiable: misma interfaz ( `.invoke` ) para OpenAI/Gemini → reduce *vendor lock-in*.

8 · LCEL (LangChain Expression Language) — cadenas modernas

Concepto	Cómo se hace con LCEL	Ventaja
Construir cadena	<code>pipeline = plantilla   llm</code>	Interfaz estándar y consistente
Invocar	<code>result = pipeline.invoke({"var": "valor"})</code>	Interfaz unificada con <code>.invoke()</code> (igual que un LLM).
Paso de variables	Diccionario con claves = <code>input_variables</code>	Más claro que kwargs dinámicos.
Salida	<code>result.content</code> (objeto <code>BaseMessage</code> )	Estandariza la respuesta con LLMs.
Sustituye a	<code>LLMChain(...).run(...)</code> (deprecado)	Evita APIs distintas, warnings y ambigüedad.

Patrón mínimo

```
from langchain.prompts import PromptTemplate
from langchain_openai import ChatOpenAI

llm = ChatOpenAI(model="gpt-4o-mini", temperature=0.5)
tpl = PromptTemplate(input_variables=["nombre"],
                    template="Saluda a {nombre} como si fueras un asistente cortés.")
chain = tpl | llm
msg = chain.invoke({"nombre": "Carlos"}) # dict, no kwargs sueltos
print(msg.content)
```

9 · Tipos de mensajes y roles en LangChain

Clase	Rol	Uso típico
SystemMessage	Sistema/Desarrollador	Instrucciones de comportamiento, tono, límites. (No se muestran al usuario.)
HumanMessage	Usuario	Preguntas y peticiones del usuario
AIMessage	Asistente	Respuestas del modelo LLM
BaseMessage	Clase padre	Base de todos los mensajes

```
from langchain.schema import AIMessage, HumanMessage, SystemMessage

# Ejemplo de uso
mensajes = [
    SystemMessage(content="Eres un asistente útil"),
    HumanMessage(content="¿Cuál es la capital de Francia?"),
    AIMessage(content="La capital de Francia es París")
]
```

10 · Streamlit para aplicaciones IA

Componente	Método	Función
Configuración página	<code>st.set_page_config()</code>	Título, favicon, layout
Elementos UI	<code>st.title()</code> , <code>st.markdown()</code>	Añadir títulos y texto
Session State	<code>st.session_state</code>	Persistir datos entre ejecuciones
Chat UI	<code>st.chat_message(role)</code>	Mostrar mensajes de chat
Input	<code>st.chat_input()</code>	Campo de entrada para usuario
Ejecutar app	<code>streamlit run archivo.py</code>	Lanza servidor local

```
import streamlit as st

# Configuración inicial
st.set_page_config(
    page_title="Mi Chatbot",
    page_icon="🤖"
)

# Inicializar estado
if "mensajes" not in st.session_state:
    st.session_state.mensajes = []

# Mostrar mensaje con rol
with st.chat_message("assistant"):
    st.markdown("Hola, ¿en qué puedo ayudarte?")
```

11 · Gestión de memoria en chatbots

Estrategia	Implementación	Consideración
Session State	<code>st.session_state["mensajes"]</code>	Persiste durante la sesión del navegador
Lista de mensajes	Almacenar como lista Python	Fácil de iterar y actualizar
Filtrado de roles	<code>isinstance(msg, SystemMessage)</code>	No mostrar mensajes del sistema
Historial completo	Enviar todos los mensajes al LLM	Mantiene contexto de conversación

```
# Patrón para gestión de memoria
if "mensajes" not in st.session_state:
    st.session_state.mensajes = []

# Mostrar historial
for mensaje in st.session_state.mensajes:
    if isinstance(mensaje, SystemMessage):
        continue # No mostrar mensajes del sistema

    role = "assistant" if isinstance(mensaje, AIMessage) else "user"
    with st.chat_message(role):
        st.markdown(mensaje.content)
```

## 12 · Generación de respuestas con contexto

Estrategia	Método	Ventajas
Invoke directo	<code>llm.invoke(st.session_state.mensajes)</code>	Simple pero sin formato específico
Con plantilla	<code>prompt   llm</code> con LCEL	Control sobre comportamiento y formato
Streaming	<code>.stream()</code> en lugar de <code>.invoke()</code>	Respuesta progresiva como ChatGPT
Variables dinámicas	<code>{"mensaje": pregunta, "historial": mensajes}</code>	Contexto estructurado

```
# Con plantilla de prompt (recomendado)
prompt = PromptTemplate(
    input_variables=["historial", "mensaje"],
    template="""Eres un asistente útil y amigable.
    Historial: {historial}
    Responde a: {mensaje}
    Asistente: """"
)

cadena = prompt | llm
respuesta = cadena.invoke({
    "historial": st.session_state.mensajes,
    "mensaje": pregunta
})
```

## 13 · Respuestas en streaming

Paso	Código	Propósito
Stream chunks	<code>for chunk in chain.stream(...)</code>	Recibir fragmentos
Concatenar	<code>full_response += chunk.content</code>	Construir respuesta completa
Actualizar UI	<code>placeholder.markdown(full_response + "▏")</code>	Mostrar progreso con cursor
Finalizar	<code>placeholder.markdown(full_response)</code>	Mostrar sin cursor al terminar

```
# Patrón completo de streaming
response_placeholder = st.empty()
full_response = ""

for chunk in cadena.stream({"mensaje": pregunta,
    "historial": st.session_state.mensajes}):
    full_response += chunk.content
    response_placeholder.markdown(full_response + "▏")

# Mostrar respuesta final sin cursor
response_placeholder.markdown(full_response)

# Guardar en memoria
```

```
st.session_state.mensajes.append(  
    AIMessage(content=full_response)  
)
```

# Página de Notas del Tema

- Esta página está pensada para que puedas anotar ideas clave, dudas y reflexiones importantes sobre el tema anterior.
-



# Tema 2 – Fundamentos de LangChain: Componentes Core

## Guía de referencia

Usa esta guía como recordatorio rápido de los componentes y APIs clave. Está pensada para cualquier proyecto, no solo para los ejemplos del curso.

### 1 · Runnables (núcleo de las “cadenas”)

Concepto	¿Qué es?	¿Cuándo usarlo?	API/Operaciones clave	Notas
Runnable	Cualquier objeto invocable con <code>.invoke(input) → output</code> .	Para encapsular pasos de procesamiento (pre/post), modelos, parsers, etc.	<code>.invoke(x)</code>	Base de todo el LangChain Expression Language (LCEL).
Pipe ( <code> </code> )	Operador de composición secuencial.	Para encadenar pasos donde la salida de uno es la entrada del siguiente.	<code>r1   r2   r3</code>	
RunnableLambda	Adapta una función Python (lambda/def) a Runnable.	Para integrar lógica propia (limpieza, formateo, validaciones).	<code>RunnableLambda(fn)</code>	Permite mezclar componentes LangChain con utilidades Python propias.

#### Patrón típico

- Envuelve funciones Python con `RunnableLambda` (p. ej., formatear/transformar).
- Conecta con otros runnables (modelo, parser) usando `|`.
- Ejecuta con `.invoke(input)` y recibe el resultado del último eslabón.

#### Buenas prácticas

- Asigna nombres descriptivos a cada paso (mejor trazabilidad).
- Aísla transformaciones de datos en funciones pequeñas y puras.
- Piensa en tipos de entrada/salida por paso (strings, dicts, listas).

### 2 · Paralelización con RunnableParallel

Objetivo	¿Cuándo usarlo?	¿Qué hace?	API/Operaciones clave	Salida
Ejecutar ramas independientes a la vez	Cuando varios pasos consumen el mismo input y no dependen entre sí (p. ej., resumir y analizar sentimiento sobre el mismo texto preprocesado).	“Abre” el flujo en ramas concurrentes y devuelve todos los resultados juntos.	<code>RunnableParallel({key: runnable, ...})</code>	dict con las claves definidas (p. ej., <code>{"resumen": "...", "sentiment_data": {...}}</code> ).

#### Cómo encaja en la cadena

`preprocesamiento | RunnableParallel({ "resumen": summaryBranch, "sentiment_data": sentimentBranch }) | mergeResults`

#### Notas y cuidados

- El input de `RunnableParallel` se reutiliza para cada rama.
- Las ramas no garantizan orden; diseña tu merge para claves, no posiciones.
- Úsalo para ahorrar tiempo cuando el modelo/recursos lo permitan.

3 · Procesamiento por lotes con `.batch(...)`

Problema	Solución	Beneficios	API/Operaciones clave	Salida
Procesar muchas entradas (reseñas, consultas, etc.) de forma eficiente.	Ejecuta la <b>misma cadena</b> sobre una <b>lista</b> de inputs simultáneamente.	Mejor uso de recursos; menos sobrecarga; <b>gestión de errores por elemento</b> (si uno falla, continúan los demás); código más simple.	<code>chain.batch([in1, in2, in3, ...])</code>	<code>list</code> de resultados en el <b>mismo orden</b> de entrada.

Cuándo preferir `.batch` sobre un bucle con `.invoke`

- Siempre que tengas 2+ elementos a procesar y no requieras control personalizado paso a paso.
- Especialmente útil con modelos remotos: reduce latencias acumuladas.

Consejos

- Mantén cada unidad de trabajo “atómica” (una reseña, una consulta, etc.).
- Diseña la cadena para devolver objetos estructurados (p. ej., dict/JSON) y facilitar el post-proceso.

4 · Prompt Templates y Prompt Engineering

Concepto	¿Qué es?	¿Cuándo usarlo?	API/Operaciones clave	Buenas prácticas
<code>PromptTemplate</code>	Plantilla de texto con partes fijas (instrucciones/contexto) y variables (placeholders).	Cuando quieras <b>separar</b> instrucciones del contenido dinámico (preguntas, datos).	<code>PromptTemplate(template=..., input_variables=[])</code> · <code>prompt.format(**vars)</code>	Define la <b>plantilla</b> en una variable aparte y luego crea el objeto. Nombra claramente cada <code>input_variable</code> .
Validación local	Probar cómo se “rellena” el prompt sin llamar al LLM.	Antes de integrar en la cadena, para verificar posiciones y sustituciones.	<code>prompt.format(...)</code>	Ahorra tiempo y errores; ideal cuando mezclas mucho texto e inputs.
Prompt engineering	Disciplina para <b>diseñar</b> prompts efectivos (claridad, rol, formato de respuesta, restricciones).	Siempre: el prompt condiciona por completo el comportamiento del LLM.	—	Itera, testea y versiona plantillas; sé explícito (rol, estilo, longitud, formato de salida).

Patrón recomendado

1. Declara `template` (texto con placeholders).
2. Crea `PromptTemplate(template=template, input_variables=[])`.
3. Valida con `prompt.format(...)`.
4. Úsalo dentro de tu cadena (con el modelo, parsers, etc.).

5 · ChatPromptTemplate (plantillas para chats)

Concepto	¿Qué es?	¿Cuándo usarlo?	API/Operaciones clave	Notas
ChatPromptTemplate	Plantilla de lista de mensajes (con roles) para LLMs conversacionales.	Cuando separas instrucciones de sistema y contenido del usuario (y opcionalmente del asistente).	ChatPromptTemplate.from_messages([...]) · prompt.format_messages(**vars)	Es la opción por defecto para flujos conversacionales.
Roles	system , human (o user ), ai (o assistant ).	Definir comportamiento, entrada del usuario y referencias de respuestas previas.	Mensajes como tuplas: ("system", "...") , ("human", "{texto}")	Compatible con modelos tipo chat; evita "concatenar" texto a mano.
Variables	Placeholders en cualquier mensaje.	Para inyectar datos dinámicos en mensajes concretos.	{nombre_variable}	Valídalo con format_messages antes de invocar al modelo.

Patrón recomendado

```
prompt = ChatPromptTemplate.from_messages([
    ("system", "Actúa como traductor ES→EN, sé preciso."),
    ("human", "{texto}")
])
messages = prompt.format_messages(texto="Hola, ¿cómo estás?")
# En cadena LCEL:
# chain = prompt | chat_model | StrOutputParser()
```

Consejos

- Mantén el mensaje de sistema claro (rol, objetivo, estilo, límites).
- Usa variables con nombres expresivos ( {consulta} , {contexto} , {idioma\_destino} ).

6 · MessagesPlaceholder (inyectar historiales y listas de mensajes)

Problema	Solución	¿Qué permite?	API/Operaciones clave	Entrada esperada
Necesitas insertar múltiples mensajes (historial, few-shots) en la plantilla.	MessagesPlaceholder dentro de ChatPromptTemplate .	Reutilizar un listado de mensajes ya tipados ( HumanMessage , AIMessage , etc.).	MessagesPlaceholder(variable_name="historial") · prompt.format_messages(historial=[...], pregunta_actual="...")	Lista de BaseMessage (no un string concatenado).

Patrón típico con contexto + nueva pregunta

```
prompt = ChatPromptTemplate.from_messages([
    ("system", "Eres un asistente útil y mantienes el contexto."),
    MessagesPlaceholder(variable_name="historial"),
    ("human", "{pregunta_actual}")
])
# luego:
messages = prompt.format_messages(
    historial=historial_conversacion, # lista de mensajes
    pregunta_actual="¿Puedes detallar la arquitectura?"
)
```

Cuándo usarlo

- Chats con memoria (persistir turnos anteriores).

- In-Context Learning / Few-shot (inyectar ejemplos como mensajes).

Buenas prácticas

- Controla el límite de tokens recortando historial (ventana deslizante, resúmenes).
- Etiqueta claramente los mensajes de ejemplo (quién pregunta/responde).

7 · Plantillas por rol reutilizables

Concepto	¿Qué es?	¿Cuándo usarlo?	API/Operaciones clave	Ventaja
SystemMessagePromptTemplate	Plantilla solo para rol de sistema.	Reutilizar instrucciones de sistema con variables (rol, tono, normas).	SystemMessagePromptTemplate.from_template("... {tono}")	Modularidad; evita duplicar texto.
HumanMessagePromptTemplate	Plantilla solo para rol humano/usuario.	Normalizar cómo formateas la entrada de usuario en distintos módulos.	HumanMessagePromptTemplate.from_template("Mi pregunta sobre {tema} es: {pregunta}")	Homogeneiza la UX del prompt.
Composición	Combinar varias plantillas por rol en un prompt de chat.	Cuando distintos componentes aportan mensajes predefinidos.	ChatPromptTemplate.from_messages([system_t, human_t])	Jerarquía y reutilización limpias.

Patrón recomendado (composición)

```
system_t = SystemMessagePromptTemplate.from_template(
    "Eres un {rol} especializado en {especialidad}. Responde con tono {tono}."
)
human_t = HumanMessagePromptTemplate.from_template(
    "Mi pregunta sobre {tema} es: {pregunta}"
)
chat_prompt = ChatPromptTemplate.from_messages([system_t, human_t])
messages = chat_prompt.format_messages(
    rol="nutricionista", especialidad="dietas veganas",
    tono="profesional pero accesible",
    tema="proteínas vegetales",
    pregunta="¿Mejores fuentes para un atleta?"
)
```

Tips

- Centraliza tus plantillas por rol en módulos reutilizables.
- Versiona cambios de plantillas (p. ej., system\_v1 , system\_v2 ) para A/B testing.

8 · Output Parsers (estructurar la salida del LLM)

Objetivo	¿Por qué es clave?	Parsers comunes	Encaje en LCEL	Salida
Forzar/convertir la salida a un formato útil.	La mayoría de apps no consumen texto libre; necesitan JSON/CSV/objetos.	JsonOutputParser , PydanticOutputParser , StrOutputParser , CsvOutputParser	Se conectan al final de la cadena: prompt   model   parser	Tipos nativos (dict/list), strings o instancias Pydantic.

JSON “estricto”

- Usa JsonOutputParser() cuando esperes un JSON válido.
- Añade instrucciones de formato al prompt (si el parser las expone) o valida post-salida.

Modelos de datos con `PydanticOutputParser`

- Define un **modelo** Pydantic con campos y tipos.
- El parser **valida y castea** la salida al modelo (lanza error si no encaja).
- Útil para **contratos** entre componentes (APIs internas, pipelines).

Patrón con Pydantic

```
from pydantic import BaseModel, Field
from langchain_core.output_parsers import PydanticOutputParser

class Analisis(BaseModel):
    resumen: str = Field(..., description="Resumen en una frase")
    sentimiento: str = Field(..., description="positivo|negativo|neutro")
    razon: str

parser = PydanticOutputParser(pydantic_object=Analisis)

# (Opcional) Instrucciones de formato para el prompt:
# format_instructions = parser.get_format_instructions()
# ...inclúyelas en el mensaje de sistema/usuario.

chain = prompt | chat_model | parser
resultado: Analisis = chain.invoke(entrada)
```

Buenas prácticas

- Siempre termina tu cadena con un parser (aunque sea `StrOutputParser` ).
- Si el parser ofrece `get_format_instructions()` , **inyéctalas** en el prompt.
- Maneja **errores de parseo** (reintentos con `with_structured_output / fallbacks`).
- Mantén los esquemas **mínimos y estables** (evita campos ambiguos).

● 9 · Pydantic (modelado y validación de datos)

Concepto	¿Para qué sirve?	¿Cómo se usa?	Ventajas
<code>BaseModel</code>	Definir <b>esquemas</b> de datos con tipos y validación.	<code>class Usuario(BaseModel): id:int; nombre:str; activo:bool=True</code>	Convierte tipos (p. ej. <code>"123"→123</code> ), valores por defecto, validación.
<code>Field</code>	Añadir metadatos/validaciones.	<code>edad:int = Field(..., ge=0, le=120, description="Años")</code>	Reglas de rango, descripciones para guiar al LLM.
Serialización	Exportar modelos a formatos comunes.	<code>.model_dump()</code> , <code>.model_dump_json()</code>	Integración sencilla con APIs/BDs.

Patrón básico

```
from pydantic import BaseModel, Field

class Usuario(BaseModel):
    id: int
    nombre: str
    activo: bool = True

u = Usuario(id="123", nombre="Ana") # convierte tipos
json_str = u.model_dump_json()
```

Buenas prácticas

- Típa **todo** (incluye listas, dicts y submodelos).
- Usa `description` en `Field` para que los LLMs entiendan qué devolver.
- Restringe con `ge` , `le` , `min_length` , etc., para salidas más fiables.

10 · Salidas estructuradas con LangChain + Pydantic

Objetivo	¿Qué hace?	API clave	Dónde encaja
Forzar al LLM a devolver objetos tipados	El modelo responde en el esquema Pydantic indicado.	<code>chat_model.with_structured_output(MyModel) → .invoke(...)</code>	Al final de la cadena LCEL: <code>prompt   chat   (structured)</code> .

Patrón recomendado

```
from pydantic import BaseModel, Field

class Analisis(BaseModel):
    resumen: str = Field(..., description="Resumen en una frase")
    sentimiento: str = Field(..., description="positivo|neutro|negativo")

chat = ChatOpenAI(model="gpt-4o-mini", temperature=0.6)
structured = chat.with_structured_output(Analisis)

resultado: Analisis = structured.invoke("Analiza: ...")
print(resultado.sentimiento)           # acceso tipado
print(resultado.model_dump_json())     # JSON válido
```

Consejos

- Incluye descripciones en todos los campos; el LLM las usa como contrato.
- Si necesitas números/arrays anidados, tipa explícitamente ( `int` , `list[str]` , submodelos).
- Maneja errores de parseo (reintentos/fallback a `StrOutputParser` ).

11 · Arquitectura de proyectos LLM (estructura sugerida)

Carpeta	Contenido típico	Por qué
<code>models/</code>	Modelos Pydantic (esquemas de salida/entrada).	Contratos estables entre módulos.
<code>prompts/</code>	<code>ChatPromptTemplate</code> , plantillas por rol, few-shots.	Reutilización y versionado de prompts.
<code>services/</code>	Utilidades: E/S de ficheros, parsers, conectores API.	Aísla dependencias externas.
<code>ui/</code>	Interfaz (CLI/Web: Streamlit/Gradio/FastAPI).	Separa presentación de lógica.

Buenas prácticas

- Empieza por el modelo de salida (qué esperas del LLM) y diseña hacia atrás.
- Versiona prompts ( `system_v1` , `system_v2` ) y testea con datos sintéticos.
- Mantén cadenas LCEL pequeñas y componibles.

12 · Diseñar esquemas Pydantic eficaces para LLMs

Elemento	Recomendación	Ejemplo
Campos string	Define propósito en <code>description</code> .	<code>nombre:str = Field(...,"Nombre completo extraído del documento")</code>
Números	Restringe rango y unidad.	<code>ajuste:int = Field(..., ge=0, le=100, description="0-100")</code>
Listas	Limita tamaño esperado (descripción) y tipo interno.	<code>habilidades:list[str] = Field(...,"Top 5-7 habilidades")</code>
Submodelos	Representa bloques (experiencia, educación, etc.).	<code>class Experiencia(BaseModel): periodo:str; rol:str</code>

Plantilla general

```
class Resultado(BaseModel):
    titulo: str = Field(..., description="Título conciso")
    items: list[str] = Field(..., description="Lista de puntos clave (3-5)")
    score: int = Field(..., ge=0, le=100, description="Puntuación global")
```

13 · Ingesta de PDFs: extracción de texto (Python)

Tarea	Herramientas	API/esqueleto	Tips
Leer PDF en memoria	PyPDF2 , io.BytesIO	PdfReader(BytesIO(pdf_bytes))	Acepta bytes desde UI/web.
Extraer texto por página	page.extract_text()	Itera for i, page in enumerate(reader.pages): ...	Añade separadores entre páginas.
Limpieza	str.strip() / normalización	Quita líneas en blanco y espacios extra.	Ahorra tokens y coste.
Errores	try/except	Devuelve mensaje claro en errores.	PDFs escaneados → OCR externo.

Patrón básico

```
from io import BytesIO
from PyPDF2 import PdfReader

def extraer_texto_pdf(pdf_bytes: bytes) -> str:
    try:
        reader = PdfReader(BytesIO(pdf_bytes))
        partes = []
        for i, page in enumerate(reader.pages, start=1):
            txt = (page.extract_text() or "").strip()
            if txt:
                partes.append(f"\n--- Página {i} ---\n{txt}\n")
        texto = "".join(partes).strip()
        if not texto:
            return "ERROR: PDF vacío o solo imágenes."
        return texto
    except Exception as e:
        return f"ERROR: al procesar el PDF: {e}"
```

Buenas prácticas

- Para PDFs escaneados usa OCR (p. ej., Tesseract/DocTR) antes de pasar al LLM.
- Recorta texto si supera el límite de tokens (resúmenes por página o chunking).
- Normaliza codificación, bullets y tablas (considera parsers específicos si necesitas estructura).

# Página de Notas del Tema

- Esta página está pensada para que puedas anotar ideas clave, dudas y reflexiones importantes sobre el tema anterior.
-



# Tema 3 – RAG y LangChain: Cargando y recuperando datos del mundo real

## Guía de referencia

Usa esta guía como recordatorio rápido de componentes, opciones y "pitfalls". Está pensada para cualquier proyecto con LangChain, no solo para el ejemplo del curso.

### 1 · LangChain Community (integraciones externas)

Qué es: Paquete separado que concentra *integraciones de terceros* (loaders, stores, conectores a APIs y BBDD).  
Relación con `langchain-core` : Core = abstracciones y protocolos; Community = *implementaciones* concretas.

#### Instalación básica

- Entorno: activa tu venv.
- Paquete: `pip install langchain-community`
- Imports típicos: `from langchain_community.document_loaders import ...`

#### Cuándo usarlo

- Cargar/procesar documentos (PDF, web, GDrive, Notion, Slack...).
- Conectarte a bases vectoriales (FAISS, Chroma, Pinecone, Milvus, Qdrant...).
- Integrarte con APIs de terceros (Google, etc.).

#### Patrón de dependencias

- Muchas integraciones requieren librerías extra (p. ej., `pypdf` , `beautifulsoup4` , SDKs de proveedores).
- Error típico: `ModuleNotFoundError` → instala la dependencia indicada por el mensaje.

#### Buenas prácticas

- Fijar versiones en `requirements.txt` .
- Aislar secretos en variables de entorno/gestores (no subir `credentials.json` a VCS).
- Crear un módulo `adapters/` donde centralizas loaders/conectores para cambiar proveedor sin tocar el resto del código.

### 2 · Document Loaders (cargar datos en objetos `Document` )

Qué hacen: Normalizan entradas externas a objetos `Document` con:

- `page_content` : texto (o contenido) del recurso.
- `metadata` : diccionario con contexto (ruta, url, título, autor, página, timestamps, etc.).

#### Flujo típico

1. Instanciar loader → 2) `load()` (o `lazy_load()` ) → 3) lista/iterable de `Document` → 4) procesar/splitear/embeddings.

#### Loaders comunes y notas

Loader	Uso típico	Dependencias frecuentes	Notas clave
<code>PyPDFLoader</code>	Leer PDFs, 1 <code>Document</code> por página	<code>pypdf</code>	Devuelve <i>lista</i> (páginas). <code>metadata</code> incluye número de página.
<code>WebBaseLoader</code>	Extraer HTML → texto	<code>beautifulsoup4</code>	Webs dinámicas pueden requerir alternativas (Selenium/Playwright).
<code>CSVLoader</code> , <code>Unstructured*</code>	Archivos ofimáticos variados	<code>unstructured</code> , <code>pandas</code>	<code>unstructured</code> maneja docx, pptx, etc.
<code>NotionDBLoader</code> , <code>SlackDirectoryLoader</code> , <code>GmailLoader</code> , etc.	SaaS	SDK/credenciales del proveedor	Requiere OAuth/keys y scopes adecuados.
Directorios ( <code>DirectoryLoader</code> )	Carpeta con múltiples ficheros	según subloader	Combínalo con <i>glob</i> y subloaders específicos.

#### API de uso

```
from langchain_community.document_loaders import PyPDFLoader, WebBaseLoader

docs_pdf = PyPDFLoader("ruta/al.pdf").load() # -> [Document, Document, ...]
docs_web = WebBaseLoader("https://ejemplo.com").load()
```

Patrones útiles

- Iterar páginas de PDF y preservar contexto: añade claves propias a `metadata` (p. ej., `doc_id` , `section` ).
- `lazy_load()` para streams grandes: procesa en pipeline sin cargar todo en memoria.
- Normalización: crea una función que homogenice `metadata` (e.g., `source` , `type` , `uri` ) para *search logs* y *retrieval filters*.

Problemas comunes

- Contenido dinámico no aparece con `WebBaseLoader` → usa un loader con navegador headless.
- Codificaciones raras → asegúrate de `encoding` correcto o limpia con `chardet` / `ftfy` .

3 · Text Splitters (fragmentar para ventanas de contexto y eficiencia)

Problema que resuelven: Los LLMs tienen *context window* limitada (p. ej., ~128k tokens). Documentos largos → errores 413/429, coste alto, olvidos y peor *attention*.

Qué hacen: Dividen contenido en *chunks* (fragmentos) *semánticamente razonables*, con *solapamiento* para preservar continuidad.

Clases comunes y cuándo usarlas

Splitter	Cuándo usarlo	Claves de configuración
<code>RecursiveCharacterTextSplitter</code>	Opción por defecto "inteligente" (capítulos → párrafos → oraciones)	<code>chunk_size</code> , <code>chunk_overlap</code> , <code>separators</code>
<code>CharacterTextSplitter</code>	División rígida por separador simple	<code>separator</code> , <code>chunk_size</code> , <code>chunk_overlap</code>
<code>TokenTextSplitter</code>	Ajuste por tokens (útil si cobras/limitas por tokens)	<code>encoding_name</code> , <code>chunk_size</code> , <code>chunk_overlap</code>
<code>MarkdownHeaderTextSplitter</code>	Respetar jerarquía Markdown (#, ##, ###)	<code>headers_to_split_on</code>
<code>SentenceTransformersTokenTextSplitter</code>	Casos con embeddings específicos/tokenización ST	<code>chunk_size</code> , <code>chunk_overlap</code>

Parámetros esenciales

- `chunk_size` : tamaño del fragmento (caracteres o tokens).
- `chunk_overlap` : solape entre fragmentos (20–15% del `chunk_size` es frecuente).
- `separators` : prioridad de cortes ( `\n\n` , `\n` , `.` , `,` , `!` , `"` ...).
- `length_function` : cómo medir longitud (caracteres vs tokenizador).
- `add_start_index` : anota índice inicial en `metadata` para trazabilidad.

Patrones recomendados

- PDF → split por página y luego por texto: combina `PyPDFLoader` con `RecursiveCharacterTextSplitter` .
- Tamaños orientativos (ajusta según caso/RAG store):
  - *Resúmenes/LLM direct query*: 1–2k tokens, `overlap` 10–20%.
  - *Embeddings para RAG*: 300–800 tokens, `overlap` 10–15%.
- Preserva contexto en `metadata` : `source` , `page` , `chunk` , `start_index` , `section` .
- Evalúa calidad de chunks: revisa cortes en tablas/listas; quizá usa splitters específicos (Markdown/HTML).

Errores y costes

- Prompt demasiado largo → 413/429 o *rate limit*: reduce `chunk_size` , filtra/recupera top-k antes de invocar el LLM.
- Coste de inferencia alto → *pre-summarization* de chunks + *map-reduce* o *refine*.

Mini-pipeline canónico

```
from langchain_community.document_loaders import PyPDFLoader
from langchain.text_splitter import RecursiveCharacterTextSplitter

docs = PyPDFLoader("doc.pdf").load() # 1 Document por página
splitter = RecursiveCharacterTextSplitter(chunk_size=1000, chunk_overlap=150)
chunks = splitter.split_documents(docs) # mantiene metadata por chunk
```

● 4 · División “inteligente” de información con RecursiveCharacterTextSplitter

Objetivo: partir documentos largos en *chunks* manejables para LLMs y para indexación (RAG), preservando coherencia.

Clase principal

- `from langchain.text_splitter import RecursiveCharacterTextSplitter`

Parámetros clave

Parámetro	Qué controla	Guías prácticas
<code>chunk_size</code>	Longitud aproximada del fragmento (caracteres por defecto)	800–1200 para embeddings; 2k–4k para prompts directos; ajusta a la <i>context window</i> .
<code>chunk_overlap</code>	Solape entre fragmentos	10–20% del tamaño del chunk (p. ej., 100–200 sobre 1k).
<code>separators</code>	Prioridad de corte (capítulo → párrafo → oración → palabra)	Déjalo por defecto o personalízalo si cortas tablas/markdown.
<code>length_function</code>	Cómo medir longitud	Útil para trabajar en <b>tokens</b> (p. ej., tiktoken).
<code>add_start_index</code>	Añade posición al <code>metadata</code>	Facilita trazabilidad y highlights.

Patrones de uso

```
splitter = RecursiveCharacterTextSplitter(
    chunk_size=1000,
    chunk_overlap=150,
)
chunks = splitter.split_documents(docs) # preserva metadata por chunk
```

Por qué el *\*overlap\**

Mitiga cortes en frases/secciones; aporta contexto al siguiente chunk.

Pitfalls

- `chunk_size` excesivo → errores 413/429 y coste alto.
- Sin *overlap* → respuestas que “pierden el hilo”.
- PDFs: primero por **página** (loader) y luego por **texto** (splitter).

Patrón “map-reduce” para resúmenes

- `map` : resume cada chunk.
  - `reduce` : sintetiza todos los resúmenes en uno final (otra invocación al LLM).
- Ventaja: escalable a documentos muy grandes.

● 5 · Embeddings: conceptos esenciales

Definición: vector numérico de longitud fija que codifica el **significado** de un texto/objeto.

Idea central: textos semánticamente similares → vectores cercanos (alta similitud); diferentes → lejanos.

Propiedades útiles

- Dimensión fija (p. ej., 384, 768, 1024, 1536, 3072...).
- Distancias comunes: **cosine**, **L2**, **dot-product**.
- Dominio: texto, código, imágenes, audio (según modelo).

Buenas prácticas

- Mismo modelo para indexar y consultar.

- Versionado del modelo de embeddings: reindexa si cambias de modelo.
- Normalización (unit-length) si la base vectorial/consulta lo requiere.

#### Antipatrones

- Mezclar embeddings de proveedores/modelos distintos en el mismo índice.
- Usar dimensiones muy grandes sin necesidad → coste/latencia y mayor ruido.

## 6 · Embeddings con LangChain (OpenAI u otros)

#### Interfaz común

- `from langchain_openai import OpenAIEmbeddings` (o proveedor equivalente)
- Todos heredan de `Embeddings` → interoperables con `VectorStore`, `Retriever`, etc.

#### Uso mínimo

```
emb = OpenAIEmbeddings(model="text-embedding-3-large")
v = emb.embed_query("texto de ejemplo") # → list[float]
V = emb.embed_documents(["t1", "t2", ...]) # batch
```

#### Comparación rápida

Opción	Ventaja	Considera
OpenAI (hosted)	Calidad, mantenimiento cero	Coste por uso, dependencia externa
Sentence-Transformers (local)	Sin datos a terceros	DevOps, GPU/CPU y memoria
Cohere/Azure/Google, etc.	Alternativas con SLAs/regiones	API keys, matices de licencia

#### Cálculo de similitud (ejemplo)

- Cosine similarity (0–1): cercano a 1 → muy similares.
- Para ranking de resultados recuperados del vectorstore.

#### Consejos de producción

- Cachear embeddings (evitar recomputar).
- Limpiar/normalizar texto antes de embeddar (lower, espacios, símbolos).
- Guardar `metadata` rico (source, page, section, chunk\_id, model\_version).

## 7 · Bases de datos vectoriales (VectorStores)

Qué son: almacenes optimizados para vectores + búsqueda por similitud/aproximada (ANN).

#### Flujo de INGESTA estándar

1. Cargar documentos (loaders).
2. Splittear en chunks (splitter).
3. Embeddings de cada chunk.
4. Indexar en el vectorstore (con su `metadata`).

#### Flujo de CONSULTA estándar

1. Convertir la query a embedding.
2. Buscar top-k vecinos (cosine/L2/dot).
3. Recuperar chunks + metadata.
4. (Opcional) Rerank con modelo cross-encoder.
5. Pasar contexto "elegido" al LLM.

#### Opciones habituales

Tipo	Ejemplos	Dónde encaja
Local/open source	Chroma, FAISS, Qdrant, Milvus	Desarrollo, POCs, on-prem, control total
Gestionadas (cloud)	Pinecone, Weaviate Cloud, Qdrant Cloud	Producción escalable, menor DevOps

Parámetros típicos del índice

Parámetro	Impacto
<code>metric</code> (cosine/L2/dot)	Afecta ranking y normalización
<code>dim</code>	Debe coincidir con el modelo de embeddings
Algoritmo ANN (HNSW, IVF, PQ)	Compromiso entre precisión, memoria y velocidad
<code>top_k</code>	Nº de vecinos; 3–10 suele ser buen punto de partida

Buenas prácticas

- **Metadata filtering:** añade `source`, `doctype`, `page`, `lang`, `tags` para búsquedas filtradas.
- **Mantenimiento:** compactación/reindexado periódico si hay muchas altas/bajas.
- **Privacidad:** si hay datos sensibles, prefiere local/on-prem o cifrado.

Errores comunes

- Desalinear `dim` con el modelo → excepción al indexar/consultar.
- Cambiar el modelo de embeddings sin reindexar → *drift* y malos resultados.
- `top_k` muy alto → latencia y coste; muy bajo → omisión de contexto relevante.

8 · VectorStores con LangChain (Chroma)

Objetivo: almacenar *chunks* embebidos y recuperar por similitud.

Paquetes e imports

- Instalación: `pip install chromadb`
- Import típico:
  - Hoy: `from langchain_community.vectorstores import Chroma`
  - Aviso de migración: pronto → `pip install langchain-chroma` y `from langchain_chroma import Chroma`
- Loaders/splitter/embeddings frecuentes:  
`PyPDFDirectoryLoader`, `RecursiveCharacterTextSplitter`, `OpenAIEmbeddings`

Creación & persistencia

Paso	API	Nota
Cargar PDFs de una carpeta	<code>PyPDFDirectoryLoader(path).load()</code>	Devuelve <code>List[Document]</code>
Splitear	<code>splitter.split_documents(docs)</code>	Preserva <code>metadata</code>
Embeddings	<code>OpenAIEmbeddings(model="text-embedding-3-large")</code>	Usa el mismo modelo para index y query
Construir índice	<code>Chroma.from_documents(docs, embedding_function=emb, persist_directory="chroma_data")</code>	Crea SQLite + carpeta
Reabrir índice	<code>Chroma(persist_directory="chroma_data", embedding_function=emb)</code>	No reingesta

Búsqueda básica

```
results = vectorstore.similarity_search(query, k=3) # → List[Document]
# Alternativas: similarity_search_with_score, max_marginal_relevance_search (MMR)
```

Parámetros prácticos

- `persist_directory` : ruta del índice en disco.
- `collection_name` : agrupa dominios distintos.
- `k` : vecinos a devolver (3–8 suele ir bien).
- `mmr` (vía `max_marginal_relevance_search` ): más diversidad, menos duplicados.

Buenas prácticas

- Añade `metadata` útil: `source` , `page` , `chunk_id` , `start_index` , `doctype` , `person_names` , etc.
- Tamaños de chunk orientativos para RAG con contratos: 1000–5000 chars, `overlap` 200–1000.
- Reindexa si cambias **modelo** o **parámetros de split**.

Pitfalls

- `k` demasiado alto → latencia/ruido.
- Duplicados semánticos → usa **MMR** o Multi-Query Retriever (abajo).
- No mezcles embeddings de distinta dimensión en la misma colección.

• 9· Retrievers en LangChain (interfaz estándar)

Qué son: envoltorios de búsqueda que devuelven `List[Document]` con una **interfaz** unificada.

Creación desde un `VectorStore`

```
retriever = vectorstore.as_retriever(
    search_type="similarity",                # "similarity" | "mmr" | "similarity_score_threshold"
    search_kwargs={"k": 2, "score_threshold": 0.0} # según el tipo
)
docs = retriever.invoke(query) # también .batch(), .stream()
```

Search types y kwargs

<code>search_type</code>	Cuándo usar	<code>search_kwargs</code> típicos
<code>"similarity"</code>	Top-k más cercanos	<code>k</code>
<code>"mmr"</code>	Diversidad / menos duplicados	<code>k</code> , <code>fetch_k</code> (candidatos), <code>lambda_mult</code> (diversidad)
<code>"similarity_score_threshold"</code>	Filtrar por calidad mínima	<code>score_threshold</code> , <code>k</code>

Ventajas frente a `.similarity_search` directo

- Intercambiable (otro backend u otro tipo de búsqueda sin tocar el resto del pipeline).
- Compatible con *retrievers* avanzados de LC.

No confundir con: *Document Loaders* (ingesta). Retriever = **consulta**, Loader = **carga inicial**.

• 10 · Multi-Query Retriever (con LLM)

Qué hace (pipeline):

1. Reformula la **query** en variantes (sinónimos, subconsultas).
2. Ejecuta búsquedas por cada variante (usando tu retriever base).
3. **Fusiona y deduplica** resultados con ayuda de un LLM.

Setup mínimo

```
from langchain_openai import ChatOpenAI
from langchain.retrievers.multi_query import MultiQueryRetriever

base = vectorstore.as_retriever(search_type="similarity", search_kwargs={"k": 2})
llm = ChatOpenAI(model="gpt-4o-mini", temperature=0)
mqr = MultiQueryRetriever.from_llm(retriever=base, llm=llm)
docs = mqr.invoke(query)
```

Cuándo usarlo

- Consultas ambiguas o vocabulario heterogéneo.
- Colecciones con duplicidad alta.

Costes y control

- 1 Llamadas a LLM → **latencia/coste**.
- Limita variantes (si el impl. lo permite) y mantén `k` razonable.
- Úsalo junto a **MMR** para máxima cobertura sin ruido.

• 11 · RAG (Retrieval-Augmented Generation)

Arquitectura canónica

1. **Retrieve**: `retriever.invoke(query)` → `docs`.
2. **Construcción de prompt**: plantilla con **contexto** (chunks + metadata) + **pregunta**.
3. **Generate**: LLM responde **usando solo el contexto** (instrucción explícita).
4. **Post-proceso**: citar fuentes, verificar cobertura, formatear respuesta.

Elementos LangChain típicos

- `PromptTemplate` / `ChatPromptTemplate`
- `RunnableParallel` / `RunnableSequence`
- `StrOutputParser`
- `retriever` (puede ser Multi-Query, MMR, etc.)

Buenas prácticas de prompt

- Instrucción de **grounding**: "Responde solo con la información del CONTEXTO; si no está, di que no sabes."
- **Citas**: incluye `source/page/chunk_id`.
- Limita tokens: recorta a los `n` mejores chunks (3–6), ordenados por score/diversidad.

Mejoras comunes

- **MMR** o **reranking** (cross-encoder) tras retrieve.
- **Fusionar resúmenes** de chunk para reducir tokens.
- **Filtros por metadata** (fecha, persona, tipo de contrato).

Métricas de calidad

- *Faithfulness / Groundedness* (no alucina).
- *Context Precision/Recall* (recupera lo necesario sin ruido).
- Tiempo de respuesta y coste por consulta.

• 12 · Arquitectura base de un sistema RAG

Componentes mínimos

Componente	Rol	Claves de implementación
VectorStore (p. ej., Chroma)	Persistencia de embeddings	<code>Chroma(persist_directory, embedding_function)</code>
Embeddings	Index y query	<b>Mismo modelo</b> en index/consulta ( <code>OpenAIEmbeddings</code> u otro)
Retrievers	Recuperación de chunks	<code>vectorstore.as_retriever(search_type=..., search_kwargs=...)</code>
LLM (consultas)	Reformulación/derivación y apoyo al retrieve	Modelo económico/determinista (p. ej. <code>gpt-4o-mini</code> , <code>temperature=0</code> )
LLM (generación)	Redacción final	Modelo más capaz; <code>temperature</code> baja-media
Prompts	Instrucciones del retrieve y de la respuesta	Separar en archivo <code>prompts.py</code>
Cadenas LCEL	Orquestación	<code>RunnablePassthrough</code> , <code>PromptTemplate</code> , `

Config centralizada ( `config.py` )

- Modelos: `EMBEDDING_MODEL` , `QUERY_MODEL` , `GENERATION_MODEL` .
- Chroma: `CHROMA_DB_PATH` .

- Retriever: `SEARCH_TYPE` ( `"mmr"` , `"similarity"` , ...), `SEARCH_K` , `MMR_LAMBDA` (0–1), `MMR_FETCH_K` .

Buenas prácticas

- No “hardcodear” rutas ni modelos en código de negocio: usa `config.py` o env vars.
- Dos LLMs ≠ obligatorio, pero **recomendable** (coste vs calidad).

• 13 · Prompts del sistema

Tipos sugeridos

Prompt	Uso	Consejos
RAG main prompt	Instruye al LLM de respuestas para usar solo el <b>contexto</b> ; cita fuentes	Variables: <code>{context}</code> , <code>{question}</code> ; exige estructura y cautela (“si no está, di que no sabes”)
Multi-Query Retriever prompt	Derivar 2–4 consultas alternativas	Especializa al dominio (sinónimos legales, partes, ubicaciones, condiciones)
Relevancia por fragmento <i>(opcional)</i>	Filtrar chunks irrelevantes antes del prompt final	Útil si hay ruido en el índice
Extracción de entidades <i>(opcional)</i>	Añadir <b>metadata</b> rico (personas, direcciones, importes, fechas)	Útil tanto en <b>ingesta</b> como en <b>pre-prompt</b>

Antipatrones

- Prompt único y genérico para todo el pipeline.
- No fijar el **formato de salida** (dificulta UI/parseo).

• 14 · Cadena LCEL: orquestación de retriever → prompt → LLM → parseo

Esqueleto recomendado

```
rag_chain = (
    {
        "context": final_retriever | format_docs, # subcadena
        "question": RunnablePassthrough()        # llega en invoke(...)
    }
    | prompt # PromptTemplate con {context} y {question}
    | llm_generation # ChatOpenAI u otro
    | StrOutputParser()
)
```

Puntos finos

- `RunnablePassthrough()` : “placeholder” para inyectar la pregunta en `invoke` .
- Subcadenas: permite **preprocesar** el output del retriever antes del prompt.
- Devuelve un único **string** (o estructura Pydantic si lo defines).

Errores típicos

- Mezclar `Document` directly en prompt sin formateo → contexto confuso.
- Usar distintos modelos de embeddings entre index y query.

• 15 · Recuperación híbrida (Ensemble Retriever)

Idea: combinar varios retrievers y **agregar puntuaciones** con pesos.

Setup típico

```
base = vectorstore.as_retriever(search_type="mmr",
    search_kwargs={"k": SEARCH_K, "lambda_mult": MMR_LAMBDA, "fetch_k": MMR_FETCH_K}
)
simr = vectorstore.as_retriever(search_type="similarity",
```



```

search_kwargs={"k": SEARCH_K}
)

mqr = MultiQueryRetriever.from_llm(retriever=base, llm=llm_queries, prompt=multiquery_prompt)

final_retriever = EnsembleRetriever(
    retrievers=[mqr, simr],
    weights=[0.7, 0.3],
    similarity_threshold=0.75, # opcional: filtra resultados débiles
)

```

#### Cuándo usarlo

- Colecciones con **duplicados** y **vocabulario variado**.
- Queries ambiguas: **MultiQuery** cubre sinónimos; **MMR** aporta diversidad.

#### Parámetros clave

Parámetro	Efecto	Rango/tip
<code>weights</code>	Importancia relativa por retriever	Suma no tiene por qué ser 1; empieza con 0.7/0.3
<code>similarity_threshold</code>	Corta ruido	0.7–0.8 suele ser buen inicio
<code>SEARCH_K</code>	Vecinos finales	2–6 según longitud de chunks
<code>MMR_LAMBDA</code>	Relevancia vs diversidad	0.6–0.8 equilibrado
<code>MMR_FETCH_K</code>	Candidatos previos a MMR	20–50 típicamente

#### Pitfalls

- Pesos mal calibrados → o muy redundante o muy disperso.
- `threshold` alto → sin resultados; bajo → ruido.

# Página de Notas del Tema

- Esta página está pensada para que puedas anotar ideas clave, dudas y reflexiones importantes sobre el tema anterior.
-

# Tema 4 – LangGraph

## Guía de referencia

Usa esta guía como recordatorio rápido de los conceptos y APIs clave. Está pensada para cualquier proyecto, no solo para los ejemplos del curso.

### 1 · Qué es LangGraph (y cómo encaja con LangChain)

Pregunta	Resumen útil	Implicación práctica
¿Sustituye a LangChain?	No. <b>Complementa</b> a LangChain usando sus interfaces (Runnable, TypedDict, etc.).	Puedes combinar <b>chains LCEL</b> dentro de nodos de LangGraph.
¿Qué añade?	<b>Grafos con ciclos, ramificación dinámica (branching), estado global compartido.</b>	Modela bucles, reintentos, bifurcaciones y memoria sin “salirte” a Python puro.
¿Cuándo usarlo?	Cuando hay <b>decisiones dinámicas, iteraciones, memoria persistente, colaboración/ paralelismo</b> entre partes.	Flujos lineales simples → LCEL; flujos complejos → LangGraph (o mixto).

### 2 · Componentes fundamentales

Componente	Qué es	¿Qué hace?	Tips
<b>Estado (State)</b>	Estructura (normalmente <code>TypedDict</code> ) compartida por todo el grafo.	<b>Lee/escrbe</b> datos accesibles por cualquier nodo.	Devuelve <b>diccionarios parciales</b> ; LangGraph <b>fusiona</b> sobre el estado.
<b>Nodo (Node)</b>	Unidad de trabajo: <b>función Python</b> o <b>Runnable</b> .	Consume/actualiza parte del estado.	Mantén nodos <b>cohesivos</b> (1 responsabilidad).
<b>Arista (Edge)</b>	Conexión <b>dirigida</b> entre nodos.	Define el <b>flujo</b> (orden, bifurcaciones, ciclos).	Usa directivas <b>START</b> y <b>END</b> para entrada/salida.

### 3 · Patrón de arquitectura (paso a paso)

- Define el estado (p.ej., `class State(TypedDict): ...`).
- Crea el grafo de estado: `graph = StateGraph(State)`.
- Declara nodos (funciones o runnables) que **reciben y/o devuelven** partes del estado.
- Añade nodos: `graph.add_node("nombre", funcion_o_runnable)`.
- Conecta: `graph.add_edge(START, "nodoA"); graph.add_edge("nodoA", "nodoB"); graph.add_edge("nodoB", END)`.
- Compila: `app = graph.compile()`.
- Invoca con **estado inicial** (dict): `result = app.invoke(initial_state)`.

4 · Estado: definición, lectura y escritura

Aspecto	Cómo hacerlo	Detalle clave
Definir	<code>class State(TypedDict): campo1: str; n: int; ...</code>	Tipado guía el <b>merge</b> de salidas de nodos.
Inicializar	<code>initial = {"campo1": "...", "n": 0, ...}</code>	Solo necesitas lo <b>mínimo</b> para arrancar el flujo.
Leer	En el nodo: acceder como <code>dict ( state["campo1"] )</code> .	Piensa el estado como <b>contexto global</b> .
Actualizar	El nodo <b>retorna dict parcial</b> : <code>{"campo1": nuevo_valor}</code>	No necesitas devolver <b>todo</b> el estado.
Persistencia	El estado viaja por todo el grafo.	Útil para <b>memoria</b> entre pasos no contiguos.

Idea clave: Cada nodo devuelve un dict parcial y LangGraph lo fusiona en el estado según el esquema.

5 · Nodos: funciones y Runnables

Tipo	Ventaja	Cuándo
Función Python <code>(state: State) -&gt; dict</code>	Simplicidad y control total.	Transformaciones, utilidades, E/S liviana.
Runnable de LangChain	Composición declarativa (LCEL), reutilización.	LLM calls, chains ya definidas, tool-calls.

Buenas prácticas de nodo

- Pura si es posible (misma entrada → misma salida).
- Pequeña y legible (una responsabilidad).
- Validación: si parseas respuestas de LLM, considera **Pydantic**/parsers estructurados.

6 · Aristas, orden y control de flujo

Acción	API/Concepto	Nota
Fijar inicio/fin	<code>add_edge(START, "A")</code> , <code>add_edge("Z", END)</code>	Requisito para la ejecución.
Secuencial	<code>A → B → C</code>	Flujo típico “pipeline”.
Branching	Arista condicional basada en <b>estado</b>	Permite <b>rutas alternativas</b> según condiciones.
Ciclos/iteraciones	<code>A → B → A</code> (con condición de parada)	Implementa <b>bucles y reintentos</b> dentro del grafo.

Diferencia con LCEL: No estás obligado a “salida del nodo = entrada del siguiente”. Los nodos consultan el **estado global**.

7 · Compilación y ejecución

Paso	Qué hace	Tips
<code>compile()</code>	Congela el grafo y valida la conexiones.	Compila <b>una vez</b> , reutiliza la app.
<code>invoke(initial_state)</code>	Ejecuta el flujo con el estado inicial.	Loguea/inspecciona el <b>estado de salida</b> .
Estados parciales	Puedes añadir campos vacíos y llenarlos en nodos.	Útil cuando inputs llegan <b>más tarde</b> .

8 · Patrón típico de uso de “LLM sobre notas/datos” (generalización)

Paso	Patrón	Qué cuidar
Extracción	Nodo LLM que extrae <b>listas</b> (participantes, temas, acciones...).	Pide <b>formato explícito</b> (CSV, JSON). Considera <b>parsers</b> .
Agregación	Nodo que construye <b>documentos</b> (acta, resumen) a partir del estado.	Lee múltiples claves del estado, <b>escribe una</b> .
Normalización	Funciones auxiliares (limpieza, split, casting).	Mantén la <b>lógica fuera</b> del prompt si es determinista.
Transcripción/ETL externo	Paso de <b>ASR</b> o carga de datos como etapa previa.	Encapsula E/S fuera del grafo y <b>vuelca</b> al estado.

- Este patrón sirve para cualquier entrada textual (reuniones, entrevistas, tickets, logs) y cualquier salida estructurada (entidades, decisiones, artefactos).

9 · Estado avanzado: Annotated Types y reducers

Cuando **varios nodos** escriben la **misma clave** del estado, por defecto la **última** escritura **sobrescribe** a las anteriores. Para **acumular/combinar** resultados, define la clave como **Annotated** con una **función reductora**:

Elemento	Descripción	Ejemplos
Tipo anotado	<code>Annotated[List[str], reducer]</code>	La clave sabe <b>cómo</b> combinar nuevos valores.
Reducer (reductora)	Función que combina <b>valor actual</b> + <b>nuevo</b> → <b>combinado</b> .	<code>operator.add</code> (listas → concat; números → suma).
Casos típicos	Logs acumulados, listas de errores, scoring incremental, merges.	Logs: <code>List[str]</code> + <code>operator.add</code> → se <b>concatenan</b> .
Custom reducers	Función Python propia si lo predefinido no encaja.	Unir <b>sets</b> (evitar duplicados), <b>merge</b> de dicts por claves.

Checklist para reducers

- Elige el **tipo base** correcto (list, set, dict, int...).
- Selecciona reducer coherente (**add**, **union**, **merge**...).
- Inicializa el estado con el **neutro** del tipo ([], set(), {}, 0...).
- Devuelve desde cada nodo **solo** la clave afectada (p.ej., `{"logs": ["Paso X ok"]}` ).

10 · Buenas prácticas (que escalan)

- Diseña **primero el estado**: nombres claros, tipos útiles, qué produce/consume cada nodo.
- Nodos idempotentes** cuando sea posible; facilita reintentos y bucles.
- Bifurcaciones explícitas**: documenta condiciones y claves del estado usadas.
- Valida salidas de LLM** con **parsers/JSON** cuando vayan al estado.
- Observabilidad**: usa una clave `logs` (Annotated + reducer) para trazas.
- Modularidad**: encapsula **E/S externas** (ASR, DB, APIs) fuera y escribe al estado.
- Composición**: integra **chains LCEL** dentro de nodos cuando convenga.
- Reutilización**: compila una vez y **reusa** la app para múltiples invocaciones.

11 · Control de flujo y decisiones (branching)

Necesito...	Qué uso en LangGraph	¿Qué hace?	Notas rápidas
Bifurcar el flujo según el estado	Arista condicional (router/"routing function")	Evalúa <code>router(state)</code> y dirige a un nodo u otro.	Define un <b>router</b> puro (sin E/S externas).
Elegir nodo de destino por nombre	Mapping de rutas	El <code>router</code> devuelve una <b>clave</b> (p. ej. <code>"par"</code> , <code>"impar"</code> ) que mapea a nombres de nodos.	Mantén <b>coherencia</b> entre claves devueltas y nombres de nodo.
Entrar por el principio con decisión	<code>START → router(...)</code>	Condiciona el <b>primer salto</b> del grafo.	Útil para <b>pre-clasificar</b> casos.
Finalizar tras cada rama	<code>add_edge("&lt;rama&gt;", END)</code>	Cierra la ejecución en cualquier rama.	Puedes reconverger varias ramas antes de <code>END</code> .

Patrón (esqueleto):

- Define estado mínimo (p. ej. `number: int` , `resultado: str` ).
- Declara nodos "hoja" (p. ej., `caso_par` , `caso_impar` ) que **solo escriben** el resultado.
- Declara `router(state)` que devuelve una **clave** de ruta ( `"par"` / `"impar"` ).
- Conecta `START` con **arista condicional** (source, `router` , `{ "par": "caso_par", "impar": "caso_impar" }` ).
- Conecta cada rama a `END` (o a pasos posteriores).

Buenas prácticas

- El router **no** debe mutar estado; que sea **determinista** y pequeño.
- Prevé **ruta por defecto** (fallback) para entradas inválidas.
- Mantén los nombres de ramas **autoexplicativos** y documenta la **condición**.

12 · Aristas condicionales (routing) y decisiones

Necesito...	API/Concepto	¿Qué hace?	Tips
Decidir camino en tiempo de ejecución	<code>graph.add_conditional_edges(source_node, router_fn, routes_map)</code>	Llama a <code>router_fn(state)</code> y salta al nodo mapeado.	El <code>router_fn</code> <b>no muta</b> estado; devuelve una <b>clave</b> (nombre de ruta).
Enrutar desde el inicio	<code>START → router_fn</code> (vía <code>add_conditional_edges</code> )	Primera bifurcación con el <b>estado inicial</b> .	Nombra rutas = <b>nombres reales de nodos</b> .
Cerrar cada rama	<code>graph.add_edge("&lt;rama&gt;", END)</code>	Finaliza tras la rama.	O re-converge en un nodo común antes de <code>END</code> .

Esqueleto de router

```
def router(state: State) -> str:
    # lógica pura en base a `state`
    return "nodo_destino" # p. ej. "par" / "impar" / "escalado"
```

Errores comunes: rutas que **no** existen; router que **modifica** el estado; strings de ruta ≠ nombres de nodos.

13- Streaming y observabilidad con LangGraph

Necesito...	API / Parámetro	¿Qué obtengo?	Cuándo usarlo
Ver progreso paso a paso	<code>app.stream(initial_state, config=..., stream_mode="updates")</code>	Emite <b>actualizaciones parciales</b> del estado por nodo.	Para mostrar "qué cambió" en cada transición.
Ver estado completo tras cada nodo	<code>stream_mode="values"</code>	<b>Estado completo</b> después de cada ejecución de nodo.	Debug detallado o auditoría.
Ver tokens/mensajes LLM	<code>stream_mode="messages"</code>	<b>Tokens/mensajes</b> generados por LLMs de los nodos.	Inspeccionar prompts/respuestas.
Debug exhaustivo	<code>stream_mode="debug"</code>	Mezcla de valores, updates y mensajes.	Diagnósticos finos (más ruidoso).

Patrón básico (loop de streaming)

```
for chunk in app.stream(state_inicial, config=config, stream_mode="updates"):
    for node_name, salida in chunk.items():
        # p.ej., acumular trazas si `salida` contiene "historial"
```

Consejos de observabilidad

- Mantén una clave `historial` (Annotated + reducer) para **trazas** legibles.
- En UI (Streamlit, FastAPI websockets, etc.), **pinta** en tiempo real lo que llegue del stream.
- Limita `stream_mode` a lo necesario (menos ruido = mejor UX).

Config & threading (para reanudación)

- Estructura típica: `config={"configurable": {"thread_id": ticket_id}}`.
- El `thread_id` identifica la **ejecución** y permite **reanudación** contra el mismo checkpoint.

14- Detención y reanudación (Human-in-the-Loop) desde la UI

Paso	Qué haces	API/Concepto	Detalles clave
1. Invocar grafo	Lanzas <code>stream()</code> con <code>state_inicial</code> y <code>config</code> (incluye <code>thread_id</code> ).	<code>app.stream(...)</code>	Si el grafo tiene <code>interrupt_before=["procesar_humano"]</code> , se <b>detiene</b> ahí.
2. Leer estado parcial	Al detenerse, obtienes <b>estado + config</b> (con identificadores de checkpoint/hilo).	Devuelto tras <code>stream()</code>	Guarda en tu UI (p. ej., <code>st.session_state["tickets"][ticket_id]</code> ).
3. Solicitar input humano	Muestras textarea/campo para <code>respuesta_humano</code> .	UI (Streamlit)	También puedes mostrar <b>contexto</b> recuperado por RAG para ayudar al agente.
4. Actualizar estado persistido	Inyectas la respuesta en el estado del hilo detenido.	<code>app.update_state(config, {"respuesta_humano": texto})</code>	No re-compilas; <b>actualizas</b> el estado guardado en SQLite.
5. Reanudar ejecución	Continúas el grafo desde el punto de pausa.	<code>app.stream(..., config=config)</code> o <code>app.invoke(..., config=config)</code>	Usa el mismo <code>config</code> / <code>thread_id</code> que al pausar.
6. Cerrar ciclo	Lees <code>respuesta_final</code> , marcas ticket como resuelto (IA u humano), guardas trazas.	Estado final	Muestra fuentes, confianza, tiempos, quién resolvió.

Checklist HITL

- Define pausa con `interrupt_before=["procesar_humano"]` (o `interrupt_after` según convenga).
- Persiste con `SqliteSaver` (o `MemorySaver` si no necesitas durabilidad).
- Usa `update_state(...)` antes de reanudar.
- Mantén **idempotencia**: si re-intentas, los nodos deberían tolerarlo.

# Página de Notas del Tema

- Esta página está pensada para que puedas anotar ideas clave, dudas y reflexiones importantes sobre el tema anterior.
-



# Tema 5 – Memoria y Gestión de Contexto con LangChain y LangGraph

## Guía de referencia

Usa esta guía como recordatorio rápido de mecanismos y parámetros para gestionar memoria y contexto en apps con LLM. Es general y aplicable a cualquier proyecto (no depende de los ejemplos del curso).

### 1 · Conceptos clave (coste, contexto y sesiones)

Concepto	¿Qué es?	Por qué importa
Context window	Límite de tokens que el LLM puede procesar por petición.	Define cuánto historial/caracteres puedes enviar; superar el límite da errores o truncado.
Coste por tokens	La facturación se basa en tokens de entrada/salida.	Más historial ⇒ más coste y latencia. Optimiza memoria y prompts.
Historial de conversación	Mensajes previos de usuario y asistente.	Mantiene coherencia (p. ej., recordar el nombre del usuario).
Sesiones / threads	Identificador para aislar conversaciones concurrentes.	Evita mezclar historiales de distintos usuarios/hilos.
Persistencia	Dónde se guarda la memoria (RAM o disco/BD).	RAM es volátil; disco permite reanudar conversaciones tras reinicios.

### 2 · Patrón básico: historial manual en el prompt (rudimentario)

Idea: Mantener una lista `history: List[BaseMessage]` y pasarla a la plantilla con `MessagePlaceholder("history")`.

Elemento	Cómo se usa	Ventajas	Limitaciones / riesgos
<code>MessagePlaceholder("history")</code>	Inserta una variable (lista de mensajes) en el prompt.	Control total del prompt.	Gestión manual tediosa; crece sin control; sin sesiones por defecto.
<code>HumanMessage</code> , <code>AIMessage</code>	Tipos de mensaje para construir <code>history</code> .	Estructura clara.	Hay que actualizar <code>history</code> tras cada turno.
Bucle de chat	<code>chain.invoke({"input": user_input, "history": history})</code>	Fácil de entender.	Sin aislamiento por sesión; alto coste si el historial crece.

Cuándo evitarlo: apps multiusuario, historiales largos, necesidad de persistencia o control fino de coste.

### 3 · LangChain (simple, sin persistencia): `RunnableWithMessageHistory`

Para prototipos o apps sencillas sin necesidad de guardar en disco.

Clase: `langchain_core.runnables.history.RunnableWithMessageHistory`

#### Piezas esenciales

- Store en RAM: diccionario `store: Dict[str, InMemoryChatMessageHistory]`.
- Historial por sesión: `InMemoryChatMessageHistory` (no persistente).
- Factory de historial:

```
def get_session_history(session_id: str):
    if session_id not in store:
        store[session_id] = InMemoryChatMessageHistory()
    return store[session_id]
```

- Envoltura con memoria:

```
chain_with_memory = RunnableWithMessageHistory(
    chain,
    get_session_history,
    input_messages_key="input",
    history_messages_key="history"
)
# tu Runnable base (prompt + LLM)
# cómo obtener el historial
# campo del input del usuario
# placeholder en tu prompt
```

- Config por sesión (aislamiento):

```
config = {"configurable": {"session_id": "<id-sesión>"}}
resp = chain_with_memory.invoke({"input": user_input}, config=config)
```

Pros / Contras

Pros	Contras
Sencillo; añade memoria y sesiones con poco código.	RAM volátil; sin persistencia; menos flexible para flujos complejos.
Compatible con prompts existentes.	No apto para cargas concurrentes/robustas a producción.

Checklist rápido

- Define `get_session_history`
- Envuelve tu cadena con `RunnableWithMessageHistory`
- Pasa `session_id` vía `config["configurable"]`

4 · LangGraph (recomendado): `MemorySaver` + `MessageState`

Para control del flujo y memoria integrada en grafos, con checkpoints.

Conceptos LangGraph

Elemento	Descripción	Uso típico
<code>StateGraph</code>	Grafo de nodos con estado compartido.	Orquestrar pasos (nodos) del agente.
<code>MessageState</code>	Estado predefinido con <code>messages: List[BaseMessage]</code> .	Mantener historial sin crear un estado custom.
Checkpoint	Guarda el estado tras cada nodo.	Inspección, reanudación, memoria.
<code>MemorySaver</code>	Checkpoint en RAM (no persistente).	Memoria rápida entre invocaciones del grafo.

Patrón mínimo

```
# 1) Definir grafo con estado de mensajes
workflow = StateGraph(StateSchema=MessageState)

# 2) Nodo de chat (lee historial y añade respuesta)
def chatbot_node(state: MessageState):
    sys = SystemMessage(content="Eres un asistente útil...")
    msgs = [sys] + state["messages"]
    resp = llm.invoke(msgs)
    return {"messages": [resp]} # se concatena al estado

workflow.add_node("chatbot", chatbot_node)
workflow.add_edge("START", "chatbot")

# 3) Compilar con checkpoint en RAM
memory = MemorySaver()
app = workflow.compile(checkpointer=memory)

# 4) Invocar por sesión (thread_id)
config = {"configurable": {"thread_id": "<id-sesión>"}}
```

```
result = app.invoke({"messages": [HumanMessage(content=user_input)]}, config=config)
```

– Ventajas clave

- Sesiones robustas mediante `thread_id` .
- Checkpointing automático (mejor inspección y reanudación).
- Performance: menos sobrecarga de “pegado” manual del historial.

– Considera persistencia

- `MemorySaver` es volátil. Para persistir, cambia el checkpointer (p. ej., SQLite/Redis) en lugar de `MemorySaver` .

• 5 · Ventana deslizante (control de coste): `trim_messages`

Objetivo: Enviar al LLM solo los N últimos turnos (manteniendo el historial completo en estado/checkpoint si quieres), para controlar tokens y latencia.

Función: `langchain_core.messages.trim_messages`

Parámetro	Qué hace	Notas
<code>strategy="last"</code>	Mantiene los últimos N “tokens” según contador.	Útil para ventana fija.
<code>max_tokens=N</code>	Límite usado por el contador.	Con contador “por mensaje”, $N \approx n^\circ$ de mensajes.
<code>token_counter=len</code>	Forma de contar tokens.	Con <code>len</code> sobre cada mensaje $\Rightarrow$ cada mensaje cuenta como 1.
<code>start_on="human"</code>	Asegura que el primer mensaje recortado sea humano.	Recomendable en chats.
<code>include_system=True</code>	Preserva el <code>SystemMessage</code> .	Mantiene el rol/instrucciones.

Patrón de uso dentro del nodo

```
# Preparar el trimmer (una vez)
trimmer = trim_messages(
    strategy="last",
    max_tokens=4,          # p. ej., últimos 4 mensajes
    token_counter=len,     # cuenta cada mensaje como 1
    start_on="human",
    include_system=True,
)

def chatbot_node(state: MessageState):
    sys = SystemMessage(content="Eres un asistente...")
    # Recortar SOLO para invocar el LLM
    trimmed = trimmer.invoke(state["messages"])
    msgs = [sys] + trimmed
    resp = llm.invoke(msgs)
    return {"messages": [resp]} # El estado aún guarda TODO el historial
```

Buenas prácticas

- Mantén estado completo para auditoría, pero recorta para invocar el LLM.
- Ajusta `max_tokens` (o cambia a contador de tokens real) según el modelo/plan.
- Usa junto a checkpointers persistentes si necesitas reanudar sesiones largas.

6 · ¿Qué enfoque usar?

Caso	Enfoque recomendado	Motivo
Prototipo rápido, 1 usuario, sin reanudar	<code>RunnableWithMessageHistory</code> + <code>InMemoryChatMessageHistory</code>	Simplicidad.
Chat multiusuario con control de flujo	<code>LangGraph</code> ( <code>StateGraph</code> + <code>MessageState</code> )	Orquestación + memoria integrada.
Memoria no persistente, rendimiento	<code>LangGraph</code> + <code>MemorySaver</code>	Rápido y aislado por <code>thread_id</code> .
Conversaciones largas / control de coste	Añade <code>trim_messages</code> (ventana deslizante)	Limita tokens enviados.
Reanudación tras reinicio/escala	Checkpoint <b>persistente</b> (SQLite/Redis/...)	Estado durable.

7 · Memoria persistente con LangGraph (SQLiteSaver)

Objetivo: conservar historiales entre ejecuciones, reinicios o despliegues.

Pieza	¿Qué es?	Opciones/Detalles
Checkpoint	Mecanismo de guardado tras cada nodo del grafo.	Se pasa al compilar: <code>app = graph.compile(checkpointer=...)</code> .
<code>MemorySaver</code>	Checkpoint <b>en RAM</b> (volátil).	Útil para prototipos, no persiste.
<code>SQLiteSaver</code>	Checkpoint <b>persistente</b> en disco.	Recomendado para conservar chats/hilos.
Conexión	<code>sqlite3.connect("historial.db", check_same_thread=False)</code>	<code>check_same_thread=False</code> permite acceso concurrente controlado.
Sesiones	<code>config={"configurable":{"thread_id": "&lt;id&gt;"}}</code>	Aísla historiales por chat/hilo. Reusar <code>thread_id</code> reanuda la conversación.

Patrón mínimo:

1. Estado y grafo

```
workflow = StateGraph(StateSchema=MessageState)
def chatbot_node(state):
    sys = SystemMessage("...instrucciones...")
    msgs = [sys] + state["messages"]
    resp = llm.invoke(msgs)
    return {"messages": [resp]}
workflow.add_node("chatbot", chatbot_node)
workflow.add_edge("START", "chatbot")
```

1. Persistencia

```
conn = sqlite3.connect("historial.db", check_same_thread=False)
cp = SQLiteSaver(conn)
app = workflow.compile(checkpointer=cp)
```

1. Invocación por sesión

```
config = {"configurable": {"thread_id": chat_id}}
app.invoke({"messages": [HumanMessage(content=user_input)]}, config=config)
```

Buenas prácticas

- Un `thread_id` por chat; un `user_id` (si lo modelas) para agrupar chats.
- Rotación/backup de la DB (vacuum, índices básicos).
- Si necesitas HA/escala: cambia a Postgres u otro checkpoint persistente equivalente.

8 · Memoria vectorial — conceptos y diseño

Idea: almacenar información en una BD vectorial y recuperar fragmentos relevantes por similitud semántica, en lugar de enviar todo el historial.

Decisión	Opciones	Comentarios
Ámbito	Por-chat, por-usuario (global), por-equipo	Lo normal: <b>global por usuario</b> + por-chat opcional.
Embeddings	OpenAI, Cohere, local (e5, Instructor...)	Elige por calidad/coste; homogeneidad entre index y query.
Store	Chroma, FAISS, PGVector, Milvus, Weaviate	Chroma: simple y persistente en disco local.
Schema	<code>document</code> , <code>id</code> , <code>metadata</code> (p. ej., <code>user_id</code> , <code>source</code> , <code>type</code> )	Etiqueta: <code>type</code> ∈ {profile, preference, work, location, ...}.
Retriever	<code>k</code> (topK), distancia (cosine, l2), filtros por metadatos	Empieza con <code>k∈[3,8]</code> ; filtra por <code>user_id</code> .
Fusión con memoria secuencial	Ventana deslizante + vectorial	Vectorial aporta <b>recuerdo duradero</b> ; ventana controla coste del turno.

Cuándo usarla

- Chats largos/antiguos donde no compensa arrastrar todo el hilo.
- Perfilado progresivo de usuario (gustos, nombre, ubicaciones, rol...).
- Búsquedas “no literales” (sin depender de coincidencia exacta).

9 · Caso práctico (abstracto): combinar memoria volátil por chat + memoria vectorial global

Arquitectura lógica

- LangGraph + `MemorySaver` → historial activo del chat (barato y rápido).
- Vector store (Chroma) → **memoria global** del usuario (perfil, gustos, datos duraderos).
- Nodo de chat:
  1. Lee último `HumanMessage` (la consulta).
  2. **Recupera** memorias globales relevantes ( `k` fijo, filtra por `user_id` ).
  3. **Inyecta** memorias a `SystemMessage` (o como contexto adicional).
  4. Invoca LLM con `[System]` + (ventana deslizante del chat) + [último Human] .
  5. **Actualiza**: añade `AIMessage` al estado + (opcional) **extrae** del último mensaje humano datos para **guardar** en vectorial.

Funciones típicas

- `save_memory(text, user_id, type=...)` → añade `document` , `id=uuid4()` , `metadata` .
- `search_memory(query, user_id, k=5)` → `vector_store.query` + filtros.
- **Heurística de guardado** (rápida) o **Clasificador con LLM** (mejor calidad) para decidir qué persistir:
  - Nombre: “me llamo...”
  - Ubicación: “vivo en... / soy de...”
  - Profesión: “trabajo como... / soy ...”
  - Preferencias: “me gusta... / me encanta...”

Consejos de implementación

- Mantén **comandos de depuración** (p. ej., `memorias` ) para listar lo guardado.
- Sanitiza/normaliza texto (minúsculas, trimming, deduplicación por hash).
- Añade **TTL** o políticas (importancia, frescura) si temes crecimiento descontrolado.
- Registra **fuelle** (chat\_id, timestamp) en metadatos.

# Página de Notas del Tema

- Esta página está pensada para que puedas anotar ideas clave, dudas y reflexiones importantes sobre el tema anterior.
-

# Tema 6 – Agentes de IA y herramientas externas con LangChain y LangGraph

## Guía de referencia

Usa esta guía como recordatorio rápido de conceptos, opciones y APIs clave. Está pensada para cualquier proyecto, no solo para el ejemplo del curso.

### 1 · Concepto de “herramienta” (tool) en LLM apps

**Idea central:** una *tool* es funcionalidad externa invocable por el LLM para hacer algo más allá de generar texto (p.ej., consultar una BD, ejecutar Python, buscar en la web).

Concepto	Qué es	Por qué importa
Tool	Función externa (con nombre, descripción y esquema de E/S) que el LLM <b>puede invocar</b> .	Permite que el modelo <i>decida</i> cuándo usar capacidades externas.
Tool vs. nodo de grafo	En un grafo (LangGraph) tú ordenas el flujo; con <i>tools</i> el LLM <b>decide</b> invocar.	Mayor flexibilidad y autonomía del agente.
Tool calling nativo	Modelos modernos generan <b>llamadas estructuradas</b> (JSON) a tools (no texto libre).	Integración más estándar y segura.
Argument schemas	Tipos y validaciones (normalmente vía Pydantic/type hints).	El LLM sabe <b>qué parámetros</b> enviar y de <b>qué tipo</b> .
Tool artifacts	Datos extra que <b>no</b> se devuelven al LLM pero quedan disponibles para lógica posterior.	Útiles para registrar metadatos, resultados intermedios o salidas pesadas (p.ej., coordenadas detecciones).

### 2 · Tool calling moderno vs. enfoque *legacy*

Enfoque	Cómo funcionaba	Problemas	Alternativa actual
<b>Legacy</b> (texto libre)	Instruías al LLM para que devolviera un “ <i>necesito usar X con argumentos Y</i> ” en texto.	Frágil, parseo propenso a errores.	
<b>Actual</b> (tool calling)	El LLM devuelve una <b>estructura JSON</b> : <code>tool_name</code> , <code>arguments</code> , <code>call_id</code> .	—	Definir tools con <b>esquemas</b> y dejar que el LLM las invoque.

Claves de definición para tool calling efectivo

- Nombre corto y preciso.
- Descripción orientada a uso (“qué hace” y “cuándo usarla”).
- Parámetros tipados (type hints / Pydantic): obligatorios, opcionales, formatos.
- Salida bien definida (string/objeto) y, si aplica, **artifacts** separados.

### 3 · Herramientas integradas y registro (LangChain)

Reutiliza tools existentes cuando sea posible. Ej.: Python REPL (intérprete Python) de `langchain-experimental.utilities.PythonREPL`.

Paso	Qué haces	Detalle clave
Instalar	<code>pip install langchain-experimental</code>	Módulo experimental con utilidades.
Crear instancia	<code>python_repl = PythonREPL()</code>	Inicializa el intérprete.
Registrar como tool	<code>Tool(name=..., description=..., func=python_repl.run)</code>	<code>name</code> y <code>description</code> guían al LLM para elegir esta tool.
Invocar manualmente	<code>tool.run("print(2+2)")</code>	Útil para testear la tool sin LLM.

Campos importantes al registrar ( `Tool` )

- `name` : único y descriptivo (p.ej., `python_repl` ).
- `description` : explicita, con *cuándo usarla* (p.ej., "ejecuta código Python para cálculos/lógica").
- `func` : *callable* que ejecuta la acción (p.ej., `python_repl.run` ).

Consideraciones

- El LLM elegirá esta tool si su descripción *encaja* con la tarea.
- Seguridad**: REPL puede ejecutar código arbitrario → sandboxing/guardrails en producción.

• 4 · Crear herramientas personalizadas con `@tool` (recomendado)

Ruta preferida por simplicidad y legibilidad.

Elemento	Requisito	Motivo
Decorador	<code>from langchain_core.tools import tool</code>	Convierte una función en tool.
Type hints	Obligatorios en parámetros y retorno	El LLM infiere tipos y valida entrada/salida.
Docstring	Describe <b>qué hace</b> y <b>cuándo usarla</b>	Se usa como <code>description</code> de la tool.

Esqueleto recomendado

- Firma: `@tool(name="...")` (opcionalmente sin nombre para usar el de la función)
- Parámetros tipados (p.ej., `query: str` )
- `-> str | dict | ...` retorno tipado
- Docstring clara: "Consulta la base de usuarios..."

Opciones útiles del decorador

Opción	Qué hace	Cuándo usarla
<code>name="..."</code>	Cambia el nombre por defecto (el de la función).	Nombres más amigables o compatibles con tu convención.
<code>return_direct=True</code>	El <b>resultado final</b> de la ejecución será la salida de la tool (no vuelve al LLM).	Flujos donde no quieres post-procesar en el LLM (p.ej., descarga de archivo).
<code>response_format=...</code>	Permite devolver <b>salida para el LLM + artifacts</b> separados.	Cuando necesitas registrar/encadenar metadatos pesados.
<code>args_schema=MyModel</code>	Define manualmente el esquema de entrada (Pydantic).	Validaciones avanzadas, alias, formatos específicos.
<code>infer_schema=True</code>	Deja que infiera esquema desde type hints.	Rápido prototipado si los hints son suficientes.

Inspección rápida (debug)

- `tool_instance.name` → nombre efectivo
- `tool_instance.description` → descripción efectiva



5 · Alternativa: **StructuredTool**

Más flexible en algunos casos, pero más verboso. Útil si quieres **construir la tool a partir de una función** sin decorarla.

Cómo	Ejemplo mental	Notas
<code>StructuredTool.from_function(fn=...) </code> o <code>StructuredTool(**kwargs)</code>	Envolvente explícito sobre una función ya existente.	Mantén type hints y docstring en la función base para que el esquema sea correcto.
Interfaz estándar	<code>.run(**kwargs)</code>	Igual que con <code>@tool</code> .

Cuándo preferir **StructuredTool**

- Necesitas **fabricar tools dinámicamente** (p.ej., a partir de un catálogo de funciones).
- Quieres **separar** definición de función y construcción de la tool por motivos de arquitectura.

6 · Patrón de diseño y buenas prácticas

- Diseña pocas tools, bien descritas. Evita solapamientos; las descripciones deben **diferenciar claramente** las capacidades.
- Valida inputs (type hints/Pydantic) y maneja errores con mensajes **útiles para el LLM**.
- Idempotencia y efectos colaterales controlados: si la tool escribe/borra, pide confirmaciones o usa *dry-run*.
- Observabilidad: registra `tool_name` , `arguments` , `duration` , `artifacts` y resultado.
- Seguridad:
  - Limita lo que el REPL puede ejecutar (listas blancas, timeouts, recursos).
  - Sanitiza entradas que viajan a SQL/APIs.
  - Mantén **quotas/reintentos** y *circuit breakers* para APIs externas.
- UX del agente: descripciones con *triggers* claros ("si necesitas cálculo aritmético...", "si debes consultar usuarios por ID...").
- Testing: prueba tools aisladas ( `tool.run(...)` ) y en **bucle LLM** (simulando invocaciones reales).

7 · Usar herramientas con LLMs en LangChain (tool calling básico)

Objetivo: habilitar que el modelo **decida** si llama a una tool y con qué argumentos.

Paso	Qué haces	API/Detalle
1) Crear el LLM	Instancia un chat model determinista.	<code>ChatOpenAI(model="gpt-4o-mini", temperature=0.2)</code> (o equivalente del proveedor que <b>soporte tool calling</b> ).
2) Tener tools	Pueden ser personalizadas ( <code>@tool</code> ) o integradas.	Asegura <b>type hints</b> + docstring; nombre y descripción claros.
3) Asociar tools al LLM	El LLM "sabe" qué herramientas puede usar.	<code>llm_with_tools = llm.bind_tools([mi_tool, ...])</code>
4) Invocar al LLM	Si el modelo necesita una tool, devolverá <code>tool_calls</code> (JSON estructurado).	<code>resp = llm_with_tools.invoke("...")</code> → <code>resp.tool_calls</code> (lista con <code>name</code> , <code>args</code> , <code>id</code> ).
5) Ejecutar la tool	Ejecuta tú la tool con esos <code>args</code> .	Preferible <code>tool.invoke(args_dict)</code> (acepta el JSON) — <code>tool.run(texto)</code> solo para una cadena simple.
6) Devolver al LLM (opcional)	Para cerrar la tarea con contexto.	Re-invoca el LLM <b>incluyendo historial</b> + resultado de la tool como mensaje/observación.

Reglas y trucos

- Nombres de tools** (especialmente con OpenAI): alfanumérico y **guiones bajos**, sin espacios ni símbolos raros. Si incumples, el proveedor **lanza error**.
- Descripción debe decir **qué hace** y **cuándo usarla**; el LLM la usa para decidir.
- `tool.run` vs `tool.invoke` :
  - `run(input_str)` → rápido para tests, **solo un string**; no gestiona `args` complejos.
  - `invoke({"arg1":..., ...})` → acepta el **payload estructurado** que devuelve el LLM.
- LCEL (cadenas)**: puedes encadenar `llm_with_tools` y luego extraer `tool_calls` (p. ej., con `operator.attrgetter("tool_calls")` o utilidades de LCEL) y mapear a `tool.invoke` para un flujo más compacto.

• 8 · Devolver un “artefacto” para uso interno (respuesta + artifact)

Objetivo: que la tool devuelva **dos salidas**: una para el modelo y otra para tu app (logs, metadatos, binarios ligeros, etc.).

Qué configurar	Cómo funciona	Nota
<code>response_format</code> en la tool	Indica que habrá <b>contenido para el LLM + artifact</b> para uso interno.	Implementa el <code>return</code> como <b>tupla</b> ( <code>contenido, artifact</code> ) .
Firma tipada	Usa <code>Tuple[str, int]</code> , <code>Tuple[str, dict]</code> , etc.	Importa <code>Tuple</code> de <code>typing</code> .
Ejecución	<code>tool.run(...)</code> solo devuelve el <b>contenido</b> .	Para inspeccionar artifact, usa <code>tool.invoke({...})</code> o procesa el mensaje estructurado.

Ejemplo mental

- **Contenido**: “Resultado de la consulta SQL ...” (entra al LLM).
- **Artifact**: `{"elapsed_ms": 123, "rows": 42}` (no entra al LLM; lo usas en observabilidad/post-procesos).

• 9 · Herramientas predefinidas e integraciones (LangChain)

LangChain incluye un gran catálogo de tools (búsqueda, ofimática, nubes, DBs, mensajería, etc.).

Categoría	Ejemplos	Cómo se usa
Búsqueda/Web	Bing, Google, Tavily, Wikipedia	Suele requerir <b>API key</b> y cliente. Instancia el wrapper y <b>ejecuta</b> con <code>.run</code> o <b>enchaina</b> con tu LLM.
Docs/Ofimática	Google Sheets/Docs, Microsoft 365	Autenticación OAuth o API key; define scopes y crea la tool (normalmente en <code>langchain_community.tools/...</code> ).
Datos/DBs	SQL, Mongo, Elasticsearch	Configura conexión; tools para <b>consultar</b> /resumir tablas.
Mensajería/Comms	Slack, Twilio, Gmail (borradores)	Requiere credenciales; útiles para agentes que <b>actúan</b> (enviar aviso, crear borrador).
Cálculo/Code	Python REPL, interpretes	<b>Peligroso</b> en producción: sandbox, timeouts y whitelists.

Patrón típico de uso

1. Instalar paquete(s) externo(s) necesario(s) ( `pip install ...` ).
2. Importar wrapper y/o tool pre-hecha.
3. Configurar con claves/parámetros.
4. Probar con `.run(...)` .
5. Añadir a `llm.bind_tools([...])` para que el agente/LLM pueda decidir usarla.

• 10 · Fundamentos de los agentes de IA (LangChain / LangGraph)

Definición: sistema dirigido por un LLM que **elige acciones** (tools) y su **orden** para lograr un objetivo.

Componente	Rol	Claves
LLM	“Cerebro” que decide el siguiente paso.	Debe soportar <b>tool calling</b> .
Toolkit	Conjunto de tools disponibles.	Limita y <b>describe bien</b> cada tool.
Plantilla de razonamiento	Patrón <b>ReAct</b> (Pensar→Actuar→Observar→...→Responder).	Viene predefinida en agentes de LangChain/LangGraph.
Memoria	Historial/estado para iterar correctamente.	Conversacional o específica de tarea.
Ejecutor del agente	Orquesta el bucle pensamiento/acción/observación.	En LangChain y LangGraph hay <b>agentes predefinidos</b> para no codificar el loop a mano.

Ciclo de un agente

- 1. Instrucción del usuario →
- 2. Pensamiento (LLM con prompt de agente) →
- 3. Elección de tool + args →
- 4. Ejecución de la tool →
- 5. Observación (resultado vuelve al contexto) →
- 6. Repite 2-5 hasta respuesta final.

Buenas prácticas

- Toolkit pequeño y bien diferenciado.
- Mensajes de error útiles para que el agente corrija.
- Guardrails: límites de pasos, coste, timeouts, rate limiting.
- Trazabilidad: log de `tool_calls` , duración, artifacts y resultados.

11 · Primer agente con LangChain (tool-calling + ejecución)

Objetivo: crear un agente que razone, elija tools y ejecute pasos automáticamente.

Elemento	API/Opción	¿Qué hace?	Notas clave
LLM	<code>init_chat_model("gpt-4o-mini", temperature=0)</code> o <code>ChatOpenAI(...)</code>	Modelo con <b>tool calling</b> nativo.	Usa temperatura baja para comportamiento estable.
Toolkit	p. ej. <code>GmailToolkit().get_tools()</code>	Colección de tools listas (buscar, leer, crear borradores...).	Evita especificar tool por tool; el agente elige.
Tools extra	<code>@tool(...)</code> funciones propias	Añade capacidades específicas del negocio.	Combina toolkit + tools personalizadas en una lista.
Prompt de agente	<code>ChatPromptTemplate</code> con <code>system</code> + <code>{input}</code> + <code>{agent_scratchpad}</code>	Define pasos, criterios y formato.	<code>agent_scratchpad</code> es <b>obligatorio</b> en agentes LC: almacena pensamiento/acciones/observaciones intermedias.
Agente	<code>create_tool_calling_agent(llm, tools, prompt)</code>	Prepara un agente ReAct que puede invocar tools.	Requiere LLM con tool-calling.
Ejecución	<code>AgentExecutor(agent=..., tools=..., verbose=True, handle_parsing_errors=True, max_iterations=K)</code>	Orquesta el bucle pensar→actuar→observar.	Define <b>límite de iteraciones</b> para evitar bucles/costos.
Invocación	<code>executor.invoke({"input": "..."})</code>	Corre el flujo completo y devuelve <code>output</code> .	Pasa <b>historial</b> si necesitas contexto conversacional.

Buenas prácticas

- Nombre y descripción de cada tool: deben dejar claro cuándo usarlas.
- Errores: `handle_parsing_errors=True` reintenta automáticamente pasos frágiles.
- Coste/seguridad: limita herramientas peligrosas (REPL), añade timeouts/quotas.

12 · Añadir tools personalizadas al agente (patrones útiles)

Caso: extender un toolkit con una tool propia (p. ej., responder en el mismo hilo de un sistema de mensajería/correo).

Paso	API/Opción	Detalle
Definir tool	<code>from langchain_core.tools import tool → @tool(name="create_reply_draft")</code>	Usa <b>type hints</b> y un <b>docstring</b> que documente: qué hace, <i>cuándo</i> usarla y cada argumento.
Acceso a APIs	Usa wrapper/cliente del servicio (p. ej., <code>toolkit.api_resource</code> )	Mantén la lógica de autenticación/config en un sitio.
Cuerpo del mensaje	Construye payload (headers, codificación, referencias)	Ej.: añadir <code>thread_id</code> / <code>in_reply_to</code> cuando el servicio lo requiera para <b>responder en hilo</b> .
Añadir al agente	<code>tools = toolkit.get_tools() + [create_reply_draft]</code>	El agente la podrá elegir según la descripción.

Checklist de una tool sólida

- **Args** bien tipados (obligatorios vs opcionales).
- **Errores** descriptivos (ayudan al agente a corregir).
- **Idempotencia** o `dry_run` si hay efectos (envíos/altas).
- **Artefactos** si necesitas telemetría (duración, IDs creados).

• 13 · Agentes con LangGraph (recomendado para producción)

Objetivo: mismo patrón ReAct pero con estado y orquestación más estables.

Elemento	API/Opción	¿Qué hace?	Notas
Agente ReAct	<code>from langgraph.prebuilt import create_react_agent</code>	Crea un agente listo para tools.	Sustituye a <code>create_tool_calling_agent</code> en LC.
Memoria de checkpoints	<code>from langgraph.checkpoint.memory import MemorySaver</code>	Guarda/recupera estado entre pasos/llamadas.	Cambiable por almacenes persistentes.
Prompt	parámetro <code>prompt=</code> (system prompt)	Define el comportamiento.	No necesitas <code>{agent_scratchpad}</code> manual.
Config	<code>config={"configurable": {"thread_id": "..."} }</code>	Aísla sesiones y reanuda estados.	Útil para multiusuario o sesiones largas.
Estado	Mensajes en <code>state["messages"]</code>	Entrada: <code>HumanMessage</code> (o <code>{"role": "user", "content": "..."} </code> )	La salida se lee del <b>último mensaje</b> del estado.

Diferencias vs LangChain

- Menos pegamento manual (scratchpad/ejecutor) y **mejor control de estado**.
- API cambia con más frecuencia: revisa parámetros actualizados (p. ej., `prompt=` en lugar de `state_modifier=` ).

• 14 · Sistemas multiagente con LangGraph Supervisor

Patrón: un supervisor selecciona/agrega respuestas de **agentes especializados**.

Componente	API/Opción	Función
Agentes especialistas	<code>create_react_agent(model, tools=[...], prompt="Eres especialista en X", name="...")</code>	Uno por área (búsqueda, cálculos, extracción, redacción...).
Supervisor	<code>from langgraph_supervisor import create_supervisor</code>	Recibe la petición, <b>decide</b> a quién delegar y cuándo parar.
Modelos	Puedes usar <b>modelos distintos</b> por agente/supervisor	Optimiza coste/rendimiento por tarea.
Compilación	<code>graph = create_supervisor(...); compiled = graph.compile()</code>	Genera el grafo ejecutable.
Ejecución	<code>compiled.invoke({"messages": [... ]})</code>	El estado final incluye trazas y mensajes.

## Topologías

- **Supervisor** (recomendada): 1 coordinador + N expertos.
- **Red**: agentes colaboran entre sí.
- **Jerárquica**: varios niveles (más compleja y frágil).

## Guardrails

- Límite de **pasos/iteraciones**, presupuesto, y tiempos.
  - Trazabilidad: log de `tool_calls`, latencias, errores, artifacts.
- 

## • 15 · Arquitectura de una solución multiagente “product-ready”

Capas clave (independientes del dominio):

### 1. Interfaz & API

- **API REST** para ingestión de eventos/tareas (POST JSON), **healthcheck**, y consulta de estados.
- Cliente UI (dashboard) que lea estados/alertas, trace pasos y resultados.

### 1. Orquestación

- **Supervisor + especialistas** (LangGraph).
- Colas/jobs opcionales si ingestión es alta (no bloquee al cliente).
- **Retentativas y backoff** para tools externas.

### 1. Integraciones

- **Búsqueda web** (proveedor search), **inteligencia** (p. ej., análisis/escáner), **mensajería/correo** (notificaciones).
- Normaliza **errores** y **rate limits** por proveedor.

### 1. Observabilidad

- Logs estructurados (agente, tool, args resumidos, duración, coste).
- Trazas por **thread\_id** y correlación con IDs de negocio.
- Métricas: latencia por tool, % reintentos, % éxito.

### 1. Seguridad y cumplimiento

- **Sandboxes** para ejecución de código.
- **Escapes/sanitización** (SQL, URLs, prompt injection).
- **Controles de salida** (validación de decisiones “riesgosas”).

### 1. Persistencia

- Estados de LangGraph (checkpoints), resultados, y artefactos.
- Configuración por **entorno** (dev/staging/prod).

# Página de Notas del Tema

- Esta página está pensada para que puedas anotar ideas clave, dudas y reflexiones importantes sobre el tema anterior.
-