



# Configuración Óptima de Mem0 para Agentes de Ingeniería

En entornos de ingeniería con múltiples agentes que manejan documentos técnicos (planos CAD, PDFs, papers, normas, etc.), es fundamental configurar Mem0 adecuadamente para gestionar la memoria de forma eficaz. A continuación se detallan las recomendaciones de expertos y mejores prácticas para configurar los distintos tipos de memoria de Mem0, consolidar hechos sin contradicciones, elegir modelos de *embeddings*, aprovechar la capa de grafo, organizar memorias por disciplina y asegurar persistencia en sesiones prolongadas. Se incluyen ejemplos de configuración JSON/YAML basados en documentación oficial.

## Tipos de memoria en Mem0: corto plazo, largo plazo, episódica y semántica

Mem0 implementa un **sistema de memoria jerarquizado** inspirado en la cognición humana. En un agente de ingeniería es recomendable habilitar y entender cada tipo de memoria:

- **Memoria de corto plazo (short-term o working memory)** – Almacena el contexto inmediato de la conversación, similar a recordar las últimas frases. Incluye el historial reciente de mensajes, variables temporales y foco actual de atención <sup>1</sup>. Esto mantiene la coherencia conversacional en el momento (ej.: recordar la última pregunta hecha por el usuario).
- **Memoria factual de largo plazo (factual long-term)** – Conserva hechos importantes, datos persistentes del usuario, preferencias, estilo de comunicación y contexto del dominio técnico <sup>2</sup>. En ingeniería, aquí se guardarían por ejemplo especificaciones confirmadas, configuraciones preferidas o conclusiones de informes previos.
- **Memoria episódica de largo plazo (episodic)** – Registra interacciones pasadas específicas y sus resultados <sup>3</sup>. Esto permite que el agente recuerde experiencias previas, como “*la última vez que se diseñó este componente, surgió un problema de vibración*”. Resulta útil para no repetir errores o recalcar aprendizajes de sesiones anteriores.
- **Memoria semántica de largo plazo (semantic)** – Almacena conocimiento generalizado o abstracto adquirido con el tiempo <sup>4</sup>. En un caso de ingeniería, podría capturar relaciones generales (p. ej., “*los procedimientos de soldadura suelen requerir un pre-calentamiento si el material es acero al carbono*” ) y entendimiento conceptual del dominio.

Estas memorias trabajan en conjunto. La memoria a corto plazo mantiene la conversación actual consistente, mientras que las memorias de largo plazo personalizan y adaptan el agente a través de múltiples sesiones <sup>5</sup> <sup>6</sup>. **Mem0** mueve dinámicamente información entre corto y largo plazo según su importancia, reciente uso y significado – imitando cómo los humanos consolidamos conocimiento <sup>7</sup> <sup>8</sup>. En la práctica, se recomienda aprovechar las capacidades integradas de Mem0 en este aspecto: el sistema aplica **filtrado inteligente** (prioriza datos importantes para evitar *bloat* de memoria) y **olvido dinámico** (decadencia de entradas de baja relevancia) para mantener la memoria eficiente <sup>9</sup> <sup>10</sup>. Esto significa que no todas las líneas de un plano o documento serán recordadas, solo aquellas destacadas por su relevancia técnica.

En resumen, para un agente de ingeniería conviene utilizar la memoria de corto plazo para el contexto inmediato (e.g. últimas interacciones del usuario), mientras la memoria larga plazo (factual, episódica y semántica) almacenará persistentemente la información técnica crucial, lecciones aprendidas de proyectos anteriores y conocimiento general de la disciplina. Mem0 se encargará de consolidar automáticamente los datos más relevantes en largo plazo y de mantener continuidad de contexto entre sesiones <sup>11</sup> <sup>12</sup>.

## Consolidación de hechos y control de contradicciones

Una funcionalidad clave de Mem0 es asegurar que la base de conocimientos del agente permanezca **coherente y sin contradicciones** incluso a medida que se agregan nuevos datos de documentos o conversaciones. El sistema compara cada nuevo hecho extraído con las memorias existentes y toma acciones de actualización apropiadas mediante un *LLM* interno <sup>13</sup>:

- **ADD (añadir):** Si la nueva información no existía en la memoria, se agrega como un nuevo recuerdo.
- **UPDATE (actualizar):** Si la información nueva se solapa con algo ya almacenado, Mem0 puede actualizar la entrada existente enriqueciendo o sustituyendo datos (por ejemplo, refinando un hecho con más detalle) <sup>14</sup> <sup>15</sup>. Se mantiene el mismo identificador de memoria al actualizar, consolidando el hecho en un solo registro.
- **DELETE (eliminar):** Si un nuevo dato **contradice** directamente un hecho previo, Mem0 marca el antiguo como obsoleto y lo elimina para resolver el conflicto <sup>16</sup> <sup>17</sup>. Por ejemplo, si la memoria decía "*El material es aluminio*" y luego un documento confirma "*El material es acero*", la entrada anterior sería eliminada por contradicción.
- **NOOP (ningún cambio):** Si el hecho recuperado ya está presente o no aporta novedad, no se realiza cambio alguno <sup>18</sup> <sup>19</sup>.

Esta lógica de consolidación de Mem0 garantiza consistencia factual. En escenarios de ingeniería, es **fundamental habilitar estas comprobaciones** para evitar que el agente se "confunda" con datos desactualizados o inconsistentes. La propia documentación de Mem0 enfatiza que el sistema detecta duplicados o contradicciones y actualiza la memoria en consecuencia, manteniéndola *coherente, sin redundancia y lista para futuras consultas* <sup>13</sup>. Se recomienda utilizar la configuración por defecto de Mem0 para este flujo de consolidación de hechos. Opcionalmente, Mem0 permite personalizar el *prompt* de actualización de memoria (*custom update memory prompt*) para ajustar estas reglas <sup>20</sup> <sup>21</sup>, aunque en la mayoría de casos técnicos las reglas estándar (añadir nuevos hechos, actualizar detalles, eliminar contradicciones) son adecuadas.

**Consejo:** Al integrar nuevos documentos técnicos (por ejemplo, subir un PDF de una norma), procure dividir la información en "hechos" discretos para que Mem0 los evalúe individualmente. Esto facilita que el motor de Mem0 detecte contradicciones entre piezas específicas de conocimiento y aplique *UPDATE* o *DELETE* correctamente en lugar de almacenar párrafos redundantes. Gracias a este enfoque, los agentes de ingeniería podrán confiar en que su memoria unificada siempre refleja la información más actualizada y verídica, evitando errores por datos viejos.

## Modelos de *embeddings* recomendados para información técnica

La calidad de los *embeddings* es crítica al trabajar con memoria semántica en documentos de ingeniería, ya que determina qué tan bien se recuperan fragmentos relevantes. **Mem0** es agnóstico al modelo de *embedding* y soporta múltiples proveedores (OpenAI, HuggingFace, Cohere, etc.), por lo que se puede elegir el más apropiado.

Para obtener la **máxima precisión semántica** en recuperación de textos técnicos, muchos expertos recomiendan el modelo de OpenAI `text-embedding-3-large`, uno de los últimos modelos de embedding de alta dimensionalidad (hasta ~3072 dimensiones). Evaluaciones independientes han encontrado que este modelo *destaca por encima del resto en calidad de embeddings*, logrando las mejores tasas de recuperación en pruebas preliminares <sup>22</sup>. La documentación de Mem0 muestra cómo configurarlo fácilmente como embedder por defecto, sustituyendo al modelo estándar más pequeño <sup>23</sup>. De hecho, si la prioridad es la **precisión** en entender similitudes entre descripciones de ingeniería, `text-embedding-3-large` suele ser la opción óptima. A continuación se ilustra cómo establecer este modelo en la configuración (usando la API de OpenAI):

```
config = {
    "embedder": {
        "provider": "openai",
        "config": { "model": "text-embedding-3-large" }
    }
}
memory = Memory.from_config(config)
```

<sup>23</sup>

No obstante, es importante considerar que los modelos más grandes implican **mayor costo y latencia** (embeddings de más dimensiones y consumo de tokens) <sup>24</sup>. Si el volumen de documentos es muy grande o se requieren respuestas en tiempo real, podría evaluarse un modelo de embeddings más ligero como compromiso. Por ejemplo, OpenAI ofrece `text-embedding-3-small` (dimensión ~1536) que es más rápido y económico, aunque con menor riqueza semántica.

Para entornos **on-premise o con restricciones de datos** (p.ej. información confidencial de ingeniería que no puede enviarse a una API externa), Mem0 permite usar modelos open-source locales a través de HuggingFace. Por defecto, el embedder de HuggingFace en Mem0 utiliza un modelo eficiente como `multi-qa-MiniLM-L6-cos-v1` <sup>25</sup> <sup>26</sup>, pero también es posible integrar modelos más potentes disponibles en *Transformers* (por ejemplo, modelos basados en BERT o SentenceTransformer entrenados en dominio técnico). Algunos candidatos de la comunidad incluyen **E5-large**, **BAAI-bge-large** o similares, que pueden ofrecer buena calidad sin depender de OpenAI. Mem0 facilita esta integración incluso vía un servicio de inferencia de embeddings de HuggingFace (Text Embeddings Inference, TEI) para acelerar el cálculo local <sup>27</sup> <sup>28</sup>.

**Recomendación:** Si el presupuesto y la latencia lo permiten, utilice **OpenAI** `text-embedding-3-large` para obtener embeddings robustos de contenido técnico (se ha observado que supera en rendimiento a otros modelos en tareas generales <sup>22</sup>). En caso de necesitar operar desconectado o reducir costos, considere modelos *open source* especializados en semántica (idealmente entrenados en textos técnicos) vía HuggingFace, configurándolos en Mem0 según la documentación. Recuerde alinear la dimensión del embedding en la configuración de la base vectorial (ej.: 1536 vs 3072) para evitar inconsistencias <sup>29</sup>.

## Uso de la capa de grafo de Mem0 (Neo4j/Memgraph) en ingeniería

Mem0 ofrece una capa opcional de **memoria de grafo** que puede activarse para representar conocimientos en forma de nodos y relaciones. Esto resulta especialmente potente en contextos de

ingeniería, donde a menudo existe una red compleja de interdependencias entre componentes, requisitos, cálculos y normativas. Al habilitar **Graph Memory**, el agente podrá almacenar relaciones explícitas entre piezas de información, permitiendo consultas más estructuradas y respuestas con mayor contexto <sup>30</sup>. En palabras de la documentación, esta integración combina *lo mejor de los enfoques vectoriales y de grafo*, logrando recuperaciones de información más precisas y completas gracias a la comprensión de relaciones entre datos <sup>30</sup>.

### ¿Conviene usar la memoria de grafo para casos de ingeniería?

En general, **sí** cuando el dominio técnico involucra muchas referencias cruzadas o topologías de conocimiento. Por ejemplo, en ingeniería mecánica un grafo podría enlazar un componente con su plano CAD, su listado de materiales y las normas que aplican; en ingeniería eléctrica podría relacionar esquemas de circuitos con cada componente y sus especificaciones. Esta representación permite que el agente responda preguntas complejas navegando esas conexiones (lo que se conoce como razonamiento *multi-hop*). De hecho, la extensión de Mem0 con grafo (denominada *Mem0g* en la investigación) añade un **módulo de extracción de entidades y relaciones** que convierte texto en grafos dirigidos: identifica entidades clave como nodos (p.ej. "válvula X") e infiere vínculos etiquetados como aristas ("válvula X está en el tanque Y") <sup>31</sup>. Con esto, las búsquedas pueden aprovechar consultas por subgrafos, encontrar rutas de relación entre conceptos, etc., algo muy valioso en ingeniería donde las piezas de información suelen estar altamente interrelacionadas.

Mem0 con grafo también maneja conflictos en ese nivel: un *detector de conflictos* marca nodos/aristas contradictorios, y el sistema decide si agregar, fusionar o invalidar elementos del grafo <sup>32</sup> (de forma análoga a ADD/UPDATE/DELETE en memoria vectorial). Esto ayuda a mantener **consistencia en conocimiento estructurado** – por ejemplo, si dos documentos establecen relaciones diferentes para un mismo componente, el sistema las reconciliará o señalará.

Para habilitar la memoria de grafo en Mem0 de forma **autoalojada**, se requiere configurar un **almacenamiento grafo** soportado, como **Neo4j** o **Memgraph** (ambos son bases de datos de grafos con soporte de lenguaje Cypher). Mem0 soporta Neo4j (incluyendo AuraDB) y Memgraph directamente; la instalación se realiza vía extras (`pip install "mem0ai[graph]"`) y proporcionando en la configuración las credenciales de la base grafo <sup>33</sup> <sup>34</sup>. A continuación un ejemplo de configuración básica con Neo4j:

```
config = {
    "graph_store": {
        "provider": "neo4j",
        "config": {
            "url": "neo4j+s://<SERVIDOR_NE04J>",
            "username": "neo4j",
            "password": "TU CONTRASEÑA"
        }
    }
}
memory = Memory.from_config(config)
```

<sup>34</sup>

En el caso de **Memgraph**, que es una alternativa de alto desempeño escrita en C++ compatible con Neo4j, se puede correr fácilmente vía Docker y apuntar Mem0 a `bolt://localhost:7687` con las

credenciales por defecto <sup>35</sup> <sup>36</sup>. Memgraph suele ofrecer mayor velocidad en gráficos en memoria, mientras Neo4j es muy maduro y estable; la elección depende de las necesidades, pero **cualquiera de los dos integrará bien con Mem0** (Mem0 simplemente ejecutará consultas Cypher bajo el capó).

**Conclusión:** Si sus agentes de ingeniería deben razonar sobre redes de información compleja (e.j., relaciones entre múltiples documentos, componentes o fases de un proyecto), **activar la memoria de grafo es recomendable**. Proveerá respuestas más “entendidas” en contexto, en lugar de solo por similitud de texto. Casos como preguntas del tipo “*¿qué normas afectan a los componentes que intervienen en el sistema X?*” o “*¿qué consecuencias tuvo cambiar el material Y en proyectos previos?*” se benefician enormemente de la estructura de grafo. Tenga en cuenta que habilitarla añade complejidad de infraestructura (se necesita correr Neo4j/Memgraph), por lo que si su caso de uso es más simple (p.ej., solo consultar párrafos puntuales de un PDF) puede prescindir de ella. En entornos de ingeniería colaborativa o con gran base de conocimiento, la inversión vale la pena por la **riqueza de respuestas contextuales** que se logra <sup>30</sup>.

## Estructuración de la memoria para agentes especializados por disciplina

En sistemas multiagente de ingeniería, es común tener agentes especializados (por ejemplo, un agente experto en temas mecánicos y otro en eléctricos). Mem0 ofrece mecanismos para **compartimentar y compartir memoria** entre estos agentes de forma flexible, usando identificadores de agente y sesión. Las buenas prácticas incluyen:

- **Separar el contexto por agente usando `agent_id`:** Al almacenar o buscar memorias, se puede especificar un `agent_id` para asociarlas a un agente particular. Esto permite crear memorias específicas de cada disciplina o rol <sup>37</sup>. Por ejemplo, un agente “mecánico” tendrá su propio conjunto de recuerdos técnicos separados del agente “eléctrico”. De este modo, cuando el agente mecánico consulta la memoria, recuperará preferentemente hechos de mecánica, sin interferencia de datos eléctricos irrelevantes.
- **Memoria compartida a nivel de usuario (`user_id`):** Además de lo anterior, Mem0 permite que ciertos recuerdos estén asociados solo a la identidad de usuario (sin definir un `agent_id`). Esos recuerdos serán **visibles para todos los agentes** del mismo usuario <sup>37</sup> <sup>38</sup>. En la práctica, esto sirve para información general o preferencias del usuario que aplican en cualquier disciplina. Por ejemplo, si el usuario “Carlos” siempre prefiere resultados en unidades métricas, ese dato se almacenaría con `user_id="carlos"` y sin `agent_id`, de forma que **todos** los agentes (mecánico, eléctrico, etc.) puedan recuperarlo.
- **Aislar información sensible o irrelevante entre agentes:** Siguiendo el punto anterior, utilice `agent_id` para que cada agente **ignore memorias ajenas a su dominio**. Mem0 facilita esto al permitir filtrar por agente en las búsquedas: p.ej., el agente eléctrico puede invocar `memory.search(query, user_id="carlos", agent_id="electrico")` y no verá memorias del agente mecánico <sup>39</sup>. Esto evita confusiones (un agente leyendo datos que no entiende) y resguarda información que no deba compartirse entre dominios.

Mem0 soporta además un parámetro `run_id` (identificador de sesión o ejecución) para incluso separar memorias por sesión de conversación o proyecto <sup>37</sup>. En un entorno de ingeniería, podría usarse `run_id` para distinguir contextos de proyecto: por ejemplo, `run_id="proyecto_A"` vs `"proyecto_B"`. Así, aunque el agente sea el mismo, no se mezclarán memorias de proyectos distintos salvo que se desee. Esta granularidad es muy útil para mantener **foco contextual**.

**Ejemplo práctico:** Supongamos un usuario “Ana” interactúa con dos agentes: *diseñador mecánico* y *diseñador eléctrico*. Durante la conversación, Ana menciona al agente mecánico: “Prefiero usar acero inoxidable en la estructura.” Eso se almacena con `user_id="ana", agent_id="mecanico"`. Luego en otro momento al agente eléctrico: “Recuerda que tengo suministro monofásico de 220V.” Se almacena con `user_id="ana", agent_id="electrico"`. Si el agente mecánico luego busca “*¿Qué preferencias materiales tiene el usuario?*”, Mem0 encontrará la memoria del acero inoxidable porque coincide el `agent_id` o está en compartido <sup>40</sup>, mientras que no vería la nota de 220V porque esa pertenece al otro agente. Pero si el agente mecánico pregunta “*¿Dónde vive Ana?*” y ese dato se almacenó solo con `user_id="ana"` (compartido), ambos agentes lo podrían recuperar <sup>39</sup> <sup>41</sup>. En resumen, use `agent_id` para especialización y memorias compartidas solo cuando proceda. La documentación de Mem0 resume que con esta estrategia uno puede: crear grafos de conocimiento por agente, compartir conocimientos comunes, aislar información especializada, y mantener contexto separado por sesiones <sup>42</sup>.

Adicionalmente, Mem0 (en su versión de plataforma gestionada) soporta categorías personalizadas que podrían mapearse a disciplinas <sup>43</sup>. En la versión open-source, esto requeriría manejar un campo `metadata` o `categories` manualmente al añadir memorias (como en el ejemplo `metadata={"category": "mecanica"}`) <sup>44</sup>. Esto no es estrictamente necesario si se utiliza bien `agent_id`, pero es otra capa posible de organización.

En definitiva, para agentes de ingeniería por disciplina: **estructure la memoria segmentándola por agente** pero permitiendo cierta superposición controlada para conocimientos universales del usuario. Esto maximiza la relevancia de la recuperación (cada agente obtiene principalmente datos de su campo) a la vez que aprovecha sinergias (compartiendo únicamente lo necesario). Mem0 está diseñado para soportar este multi-agente con un solo *backend* de memoria compartido <sup>45</sup> <sup>46</sup>, por lo que no hace falta instancias separadas de Mem0; basta con utilizar los identificadores de contexto adecuados en cada operación.

## Almacenamiento persistente, TTLs y manejo de sesiones prolongadas

En aplicaciones de ingeniería es común que los usuarios vuelvan recurrentemente y mantengan **sesiones largas o múltiples sesiones a lo largo del tiempo**. Por eso, es crucial que la memoria del agente sea **persistent**e y maneje adecuadamente la expiración de información obsoleta:

- **Bases de datos vectoriales persistentes:** Mem0 soporta integraciones con numerosos motores de vector DB (Qdrant, Chroma, Pinecone, Milvus, etc.) <sup>47</sup> <sup>48</sup>. Por defecto, si no se configura nada, Mem0 utilizará Qdrant localmente <sup>49</sup> <sup>50</sup>. Para producción, se recomienda explícitamente configurar un almacén vectorial persistente. Por ejemplo, con Qdrant se puede especificar `on_disk: True` para que almacene en disco (archivo) en lugar de solo en memoria RAM <sup>51</sup>. Esto garantizará que si el servidor de agentes se reinicia o cae, las memorias no se pierdan. A continuación un ejemplo de configuración de Mem0 con Qdrant local persistente:

```
{
  "vector_store": {
    "provider": "qdrant",
    "config": {
      "collection_name": "mem0",
      "host": "localhost",
    }
  }
}
```

```

        "port": 6333,
        "on_disk": true
    }
}
}

```

52 29

En entornos empresariales, puede optarse por instancias administradas (p. ej. Qdrant Cloud, Pinecone SaaS) o una base de datos SQL con extensión vectorial como PGVector. Mem0 ofrece flexibilidad para apuntar al backend deseado, solo asegúrese de proveer las credenciales y parámetros de conexión en el `config` (host, puerto, api\_key, etc.) <sup>53 54</sup>. Un vector DB bien dimensionado permitirá escalar a grandes conjuntos de documentos técnicos sin degradar la latencia de búsqueda.

- **Control de expiración (TTLs):** No toda información técnica debe guardarse para siempre. Mem0 permite asignar un **Time-To-Live** individual a los recuerdos mediante el campo `expiration_date` (fecha de expiración) <sup>55</sup>. Esto significa que podemos, por ejemplo, hacer que ciertas memorias “temporales” (quizá referentes a un proyecto en curso) se eliminen automáticamente después de X días. La API de Mem0 admite establecer esta fecha al agregar memorias; tras pasar la fecha, esas entradas ya no aparecerán en búsquedas. Esta práctica es útil para **mantener la base de conocimiento limpia**: por ejemplo, si se cargaron datos de sensores en tiempo real para diagnóstico, podría asignárseles TTL de 24 horas, mientras que conclusiones finales de un diseño se almacenan sin expiración. El uso de TTL también ayuda con cumplimiento de políticas (p.ej. eliminar datos de usuario tras cierto tiempo). En la plataforma gestionada de Mem0, las políticas de retención se pueden automatizar, pero en self-hosted es buena idea implementar un proceso (o usar las APIs `delete`, `delete_all`, etc.) para limpiar memorias expiradas si la base vectorial no las elimina por sí sola <sup>55 56</sup>.
- **Sesiones de usuario prolongadas:** Mem0 fue concebido para *persistir contexto a través de sesiones*, así que soporta naturalmente conversaciones largas que se extienden en el tiempo <sup>12</sup>. A diferencia de agentes tradicionales que “olvidan” todo al resetear la sesión, Mem0 mantiene el historial relevante del usuario almacenado y accesible. Para aprovechar esto:
  - Use un identificador de usuario consistente (`user_id`) para cada persona que interactúa con el sistema. De ese modo, aunque cierren la aplicación y vuelvan otro día, el agente recuperará sus recuerdos previos. Por ejemplo, `user_id="ana_lopez"` puede tener meses de interacción acumulada.
  - Aproveche la capacidad de **resumen y compendio** que Mem0 realiza automáticamente: en lugar de guardar cada mensaje de una sesión larga literalmente, Mem0 genera resúmenes periódicos y extrae los hechos clave <sup>57</sup>. Esto permite que incluso tras cientos de intercambios, la memoria permanezca manejable y las búsquedas sigan encontrando la información esencial (sin incluir cada detalle irrelevante).
  - Si se requiere *resetear* una sesión sin perder la historia global, puede usarse el `run_id` para iniciar un nuevo hilo de conversación aislado pero que aún comparte el `user_id`. Así, la conversación arranca limpia de corto plazo, pero el agente aún puede acceder a memorias antiguas del usuario si son pertinentes.

En pruebas de laboratorio (benchmark **LOCOMO** citado por Mem0), se evidenció que este enfoque de memoria persistente y selectiva produce mejoras significativas en la **precisión de las respuestas (+26%)** comparado con depender solo del contexto plano de las conversaciones <sup>58 59</sup>, a la vez que

reduce drásticamente la latencia y el costo de tokens porque no necesita mandar todo el historial al modelo en cada pregunta. Por ello, para sesiones prolongadas de ingeniería, **confiar en Mem0 para administrar el contexto** es más eficaz que aumentar desmesuradamente la ventana de contexto del LLM.

**Resumiendo:** utilice un **vector store persistente** para que la memoria sobreviva reinicios; establezca **TTLs** en datos que caducan (si aplica); y deje que Mem0 maneje la continuidad entre sesiones mediante su diseño de memoria a largo plazo. Esto garantiza que, aunque un proyecto de ingeniería se extienda por semanas, el agente seguirá recordando los detalles importantes de conversaciones pasadas (qué supuestos se acordaron, qué cálculos ya se validaron, etc.) sin contradecirse ni olvidar, brindando así una experiencia consistente y confiable al usuario.

## Ejemplo de configuración unificada de Mem0

A continuación se muestra un ejemplo integrando varias de las recomendaciones anteriores en un archivo de configuración (formato JSON/YAML) para un sistema multiagente de ingeniería. Esta configuración emplea el embedder de OpenAI con `text-embedding-3-large`, un vector DB Qdrant local persistente, y la capa de grafo con Neo4j, además de separar memorias por agente:

```
embedder:
  provider: openai
  config:
    model: text-embedding-3-large

vector_store:
  provider: qdrant
  config:
    collection_name: mem0
    host: localhost
    port: 6333
    on_disk: true

graph_store:
  provider: neo4j
  config:
    url: "neo4j+s://<NEO4J_AURA_URL>"
    username: neo4j
    password: "<TU_PASSWORD>"

# Opcional: instrucciones personalizadas y exclusiones
custom_instructions: "No almacenes datos de tarjetas de crédito ni PII del
usuario."
excludes: ["clave API", "contraseña", "dato confidencial"]
```

En este ejemplo:

- La sección `embedder` define el modelo de embedding a usar (OpenAI 3-large) <sup>23</sup>.
- `vector_store` apunta a un servidor Qdrant local, persistiendo la colección en disco <sup>52</sup> <sup>51</sup>.

- `graph_store` configura la conexión a Neo4j (podría ser local o AuraDB en la nube) para habilitar la memoria de grafo <sup>34</sup>.
- Se incluyen además dos parámetros opcionales: `custom_instructions` y `excludes`. Estas configuraciones (tomadas de la documentación avanzada) sirven para indicarle al LLM de Mem0 que **no almacene cierta información sensible** <sup>60</sup> – en el ejemplo, se le instruye ignorar datos como contraseñas o PII. Esto puede ser útil en entornos donde algunos documentos técnicos incluyan datos que por políticas no deban persistirse en memoria.

**Fuente:** Esta configuración se basa en la documentación oficial de Mem0 (repositorio y guías) <sup>23</sup> <sup>52</sup> <sup>34</sup>, adaptada a un caso de uso de ingeniería.

Por supuesto, los valores exactos (como host/puerto o claves API) deberán ajustarse a su entorno. Una vez cargada esta config, su aplicación iniciaría Mem0 así: `memory = Memory.from_config(config)`, y luego cada agente utilizaría métodos como `memory.add(...)` y `memory.search(...)` pasando los `user_id` / `agent_id` correspondientes.

---

**Referencias:** Las recomendaciones anteriores han sido compiladas a partir de la documentación técnica de Mem0 y fuentes expertas. Por ejemplo, la descripción de tipos de memoria proviene del blog oficial de Mem0 <sup>1</sup> <sup>3</sup>; las estrategias de consolidación de hechos y eliminación de contradicciones están basadas en las directrices de actualización de memoria de Mem0 <sup>13</sup> <sup>16</sup>; la selección de *embeddings* se apoya en evaluaciones externas donde OpenAI 3-large destaca en rendimiento <sup>22</sup>; el uso de la capa de grafo se respalda en la investigación publicada de Mem0 (Mem0<sup>g</sup>) <sup>31</sup> y en la guía de integración con Neo4j/Memgraph <sup>30</sup>; finalmente, las mejores prácticas multiagente y de persistencia se fundamentan en ejemplos oficiales <sup>37</sup> <sup>29</sup> y se alinean con la filosofía de Mem0 de memoria persistente y escalable <sup>58</sup> <sup>59</sup>. Estas fuentes verificables aseguran que la configuración propuesta sigue el **estado del arte** en gestión de memoria para agentes inteligentes en dominios técnicos. <sup>9</sup>

<sup>22</sup>

---

<sup>1</sup> <sup>2</sup> <sup>3</sup> <sup>4</sup> <sup>5</sup> <sup>6</sup> <sup>7</sup> <sup>8</sup> <sup>9</sup> <sup>10</sup> <sup>11</sup> <sup>12</sup> AI Agent Memory: What, Why and How It Works | Mem0  
<https://mem0.ai/blog/memory-in-agents-what-why-and-how>

<sup>13</sup> <sup>31</sup> <sup>32</sup> <sup>57</sup> <sup>58</sup> <sup>59</sup> AI Memory Research: 26% Accuracy Boost for LLMs | Mem0  
<https://mem0.ai/research>

<sup>14</sup> <sup>15</sup> <sup>16</sup> <sup>17</sup> <sup>18</sup> <sup>19</sup> <sup>20</sup> <sup>21</sup> Custom Update Memory Prompt - Mem0  
<https://docs.mem0.ai/open-source/features/custom-update-memory-prompt>

<sup>22</sup> <sup>24</sup> How to Choose the Best Embedding Model for Your LLM Application | by Apoorva Joshi | MongoDB | Medium  
<https://medium.com/mongodb/how-to-choose-the-best-embedding-model-for-your-lm-application-2f65fcdfa58d>

<sup>23</sup> OpenAI - Mem0  
<https://docs.mem0.ai/components/embedders/models/openai>

<sup>25</sup> <sup>26</sup> <sup>27</sup> <sup>28</sup> Hugging Face - Mem0  
<https://docs.mem0.ai/components/embedders/models/huggingface>

<sup>29</sup> <sup>44</sup> <sup>51</sup> <sup>52</sup> Qdrant - Mem0  
<https://docs.mem0.ai/components/vectordbs/dbs/qdrant>

<sup>30</sup> <sup>33</sup> <sup>34</sup> <sup>35</sup> <sup>36</sup> <sup>37</sup> <sup>38</sup> <sup>39</sup> <sup>40</sup> <sup>41</sup> <sup>42</sup> Overview - Mem0  
[https://docs.mem0.ai/open-source/graph\\_memory/overview](https://docs.mem0.ai/open-source/graph_memory/overview)

43 55 56 60 **Mem0: The Comprehensive Guide to Building AI with Persistent Memory - DEV Community**

<https://dev.to/yigit-konur/mem0-the-comprehensive-guide-to-building-ai-with-persistent-memory-fbm>

45 46 **LlamaIndex Multi-Agent Learning System - Mem0**

<https://docs.mem0.ai/examples/llamaindex-multiagent-learning-system>

47 48 49 50 **Overview - Mem0**

<https://docs.mem0.ai/components/vectordbs/overview>

53 54 **Configurations - Mem0**

<https://docs.mem0.ai/components/vectordbs/config>