

Agentes de Inteligencia Artificial

Nombre del taller: **INGENIERIA DE SOLUCIONES CON INTELIGENCIA ARTIFICIAL 002D**

Nombre de la asignatura: **INGENIERÍA DE SOLUCIONES CON INTELIGENCIA ARTIFICIAL**

Nombre del profesor: **GIOCRISRAI GODOY BONILLO**

Integrantes del grupo (nombre y apellido): [MATIAS NICOLAS CERDA REYES](#)
[JAVIER ALEXIS MUNOZ TORO](#)

Índice

ANÁLISIS DEL CASO ORGANIZACIONAL	4
Contexto	4
Stakeholders	4
Requerimientos	5
Restricciones	5
PROBLEMA A RESOLVER Y MÉTRICAS DE ÉXITO.....	6
Métricas de éxito ampliadas	6
Diseño de la solución LLM + RAG	7
Formulación de prompts.....	7
1. Prompt Inicial: Este prompt se usa al arrancar el agente. Define el rol, el estilo de respuesta y los límites. Su tono es claro y práctico, orientado a que el modelo entienda cómo debe interactuar con usuarios reales.....	7
2. Prompt de seguimiento: Este es el prompt que el agente usa cuando ya existe una conversación en curso. Su objetivo es leer el historial, evitar repeticiones y continuar la interacción sin perder coherencia.....	7
3. Prompt para activar herramientas: Este prompt sirve cuando el agente debe decidir si debe usar una herramienta (por ejemplo, leer notas, escribir una nueva, consultar conocimiento interno, actualizar un registro, etc.). El estilo está orientado a decisiones claras y acciones seguras.....	8
RAG — Pipeline.....	8
Arquitectura.....	10
1. Ingestion Layer.....	10
2. Vector DB.....	10
3. Retrieval Layer	10
4. LLM Layer	10
5. Agent Layer.....	11
6. Observabilidad	11
7. Seguridad.....	11
Justificación técnica	11
DESARROLLO DEL AGENTE FUNCIONAL	12
Integración de herramientas.....	12
Tool: Consultar base de datos de personal	12
Tool: Registrar ticket	12
Tool: Guardar nota persistente	12
Memoria y recuperación.....	13
Buffer memory.....	13
Summarization memory.....	13
Episodic memory.....	13
Estrategias de planificación.....	14
1. Entender la intención.....	14

2. Clasificar: información → acción → herramienta.	14
3. Seleccionar si corresponde usar:.....	14
○ RAG	14
○ Memoria	14
○ Tool	14
4. Reensamblar contexto.....	14
5. Generar respuesta final con orientación institucional.....	14
Orquestación.....	15
Métricas ampliada	16
Latencia	16
Precisión	16
Consistencia.....	16
Logs y análisis	17
Hallazgos:	17
Protocolos de seguridad.....	17
CONCLUSIONES + MEJORAS PROPUESTAS	18
Conclusiones	18
Mejoras y propuestas futuras	19
1. Rendimiento	19
2. Calidad de recuperación	19
3. Memoria	19
4. Arquitectura	19
5. Seguridad	19
6. Experiencia de usuario	19
ANEXOS.....	20
Capturas	20

ANÁLISIS DEL CASO ORGANIZACIONAL

Contexto

Las organizaciones modernas manejan un volumen creciente de información interna: políticas, protocolos, instructivos, normativa legal, procesos operacionales y documentación administrativa.

El acceso a esta información suele ser fragmentado — repositorios, carpetas compartidas, PDFs, correo interno — generando:

- Duplicación de preguntas entre colaboradores.
- Baja productividad del personal administrativo.
- Pérdida de trazabilidad en las consultas internas.
- Retrasos operacionales en respuestas críticas.

El proyecto nace para centralizar, automatizar y estandarizar la entrega de conocimiento institucional mediante un agente conversacional potenciado con LLM + RAG.

Es un asistente capaz de interpretar preguntas naturales, buscar información en documentos internos, recuperar memoria, ejecutar herramientas (APIs, DB, operaciones internas) y generar respuestas estructuradas con acompañamiento de trazabilidad y observabilidad.

Stakeholders

- Estudiantes
- Personal administrativo
- Dirección de TI
- Coordinación académica
- Equipo desarrollador del agente

Requerimientos

- Generación de respuestas consistentes basadas en fuentes verificadas.
- Sistema capaz de actualizar su conocimiento sin reentrenar el modelo.
- Persistencia de memoria para mejorar las interacciones.
- Logs detallados de consultas para auditoría.
- Interfaz API simple para integrarse con portales internos o dashboards.
- Bajos costos de operación y escalabilidad progresiva.

Restricciones

- El sistema debe respetar políticas de privacidad corporativa.
- Los documentos internos pueden contener información sensible, por lo tanto:
 - RAG debe ser local u on-premise.
 - No se puede enviar documentos completos al modelo.
- Infraestructura inicial de bajo costo → se priorizan componentes livianos.
- Tiempo limitado de respuesta: SLA de < 2 segundos para preguntas simples.

PROBLEMA A RESOLVER Y MÉTRICAS DE ÉXITO

Las organizaciones requieren un medio eficiente para democratizar el acceso al conocimiento interno y reducir la sobrecarga operacional. El proyecto busca reemplazar consultas repetitivas y dispersas por un asistente centralizado capaz de proporcionar:

- Respuestas inmediatas.
- Información siempre alineada a los documentos oficiales.
- Interacción natural y en lenguaje cotidiano.
- Mecanismos de control, trazabilidad y aprendizaje continuo.

Métricas de éxito ampliadas

- Precisión RAG
Objetivo $\geq 85\%$. Se evalúa comparando chunks relevantes vs chunks recuperados.
- Reducción de carga operacional
Ahorro esperado $\geq 40\%$ en consultas internas repetitivas.
- Latencia promedio
Meta ≤ 1.5 segundos por respuesta.
- Índice de satisfacción del usuario
Encuestas internas $\geq 4/5$.
- Tasa de resolución autónoma
Meta $\geq 70\%$ sin intervención humana.

- Tasa de alucinación controlada
Meta $\leq 5\%$ de respuestas sin evidencia documental.

Diseño de la solución LLM + RAG

Formulación de prompts

1. **Prompt Inicial:** Este prompt se usa al arrancar el agente. Define el rol, el estilo de respuesta y los límites. Su tono es claro y práctico, orientado a que el modelo entienda cómo debe interactuar con usuarios reales

“Eres Eva, asistente virtual de la organización. Tu función es ayudar a empleados y estudiantes entregando información confiable, basada en los documentos internos que tengas disponibles. Explica las cosas de forma sencilla, pide aclaraciones cuando algo no sea claro y evita inventar datos. Si el usuario solicita acciones específicas, guíalo paso a paso o utiliza las herramientas internas del sistema cuando corresponda.”

2. **Prompt de seguimiento:** Este es el prompt que el agente usa cuando ya existe una conversación en curso. Su objetivo es leer el historial, evitar repeticiones y continuar la interacción sin perder coherencia

“A continuación tienes el historial reciente de la conversación. Úsalo para mantener el contexto y continuar la respuesta del modo más coherente posible. Evita repetir información ya entregada y revisa si el usuario dejó alguna solicitud pendiente. Si detectas inconsistencias o falta de datos, pide una aclaración antes de avanzar.”

- 3. Prompt para activar herramientas:** Este prompt sirve cuando el agente debe decidir si debe usar una herramienta (por ejemplo, leer notas, escribir una nueva, consultar conocimiento interno, actualizar un registro, etc.). El estilo está orientado a decisiones claras y acciones seguras

“Evalúa si la petición del usuario requiere usar una herramienta del sistema.

– Si corresponde, elige la herramienta adecuada y especifica claramente los argumentos.

– Si la información puede responderse con los documentos internos, usa el contenido recuperado por RAG.

– Si no existe suficiente información, indícalo.

No inventes resultados y prioriza siempre acciones seguras y coherentes.”

RAG – Pipeline

Flujo completo

1. • Carga de documentos en `storage/docs/*.md`.
2. • Limpieza y preprocesamiento.
3. • Generación de embeddings.
4. • Almacenamiento en **ChromaDB**.
5. • Consulta semántica.
6. • Reranking por similitud.
7. • Envío de chunks al LLM.

Manejo de versiones de documentos

Se recomienda almacenar:

- *Hash MD5.*
- *Fecha de ingestión.*
- *ID de versión.*

Esto permite trazabilidad de respuestas y auditoría completa.

Arquitectura

1. Ingestion Layer

- *Parser Markdown → texto plano.*
- *Chunker.*
- *Generación de embeddings.*
- *Registro en metadata store.*

2. Vector DB

- *Metadatos de documentos.*

3. Retrieval Layer

- *Similarity search.*
- *Filtros por categoría (políticas/procedimientos).*
- *Re-ranking opcional.*

4. LLM Layer

- *Modelo principal.*
- *Motor de prompts.*
- *Control de formato de salida.*

5. Agent Layer

- *Planificador (planning.py).*
- *Memoria conversacional y episódica.*
- *Selector de herramientas.*

6. Observabilidad

- *Dashboard web ([dashboard.py](#)).*

7. Seguridad

- *Sanitización de entradas.*

- *Cifrado de claves.*

Justificación técnica

Por qué RAG y no fine-tuning

- *Se evita el riesgo de filtrar datos internos.*
- *Actualizar políticas es inmediato.*
- *Mucho más económico.*
- *Respuestas trazables y explicables.*

¿Por qué embeddings?

Permiten identificar similitud semántica entre preguntas y contenido, incluso si el usuario utiliza lenguaje informal.

DESARROLLO DEL AGENTE FUNCIONAL

Integración de herramientas

Ejemplos prácticos:

1. Web Search (DuckDuckGo)

El agente activa esta herramienta cuando la información no está en documentos internos.

2. RAG mediante ChromaDB

Consulta documentos institucionales indexados.

3. LLM vía Ollama

Se usa para generar la respuesta final en todos los modos.

Memoria y recuperación

Buffer memory

- *Conserva la última ventana de conversación (5–10 interacciones).*
- *Ideal para follow-ups.*

Episodic memory

- *Datos importantes y permanentes del usuario.*
- *Guardada en `storage/notes.json`.*

TTL recomendado

Tipo	TTL	Motivo
<i>Conversacional</i>	<i>20 minutos</i>	<i>Mantener coherencia</i>

<i>Episódica</i>	<i>Permanente</i>	<i>preferencia o datos estables</i>
<i>Herramientas</i>	<i>Según sesión</i>	<i>Evitar contaminación</i>

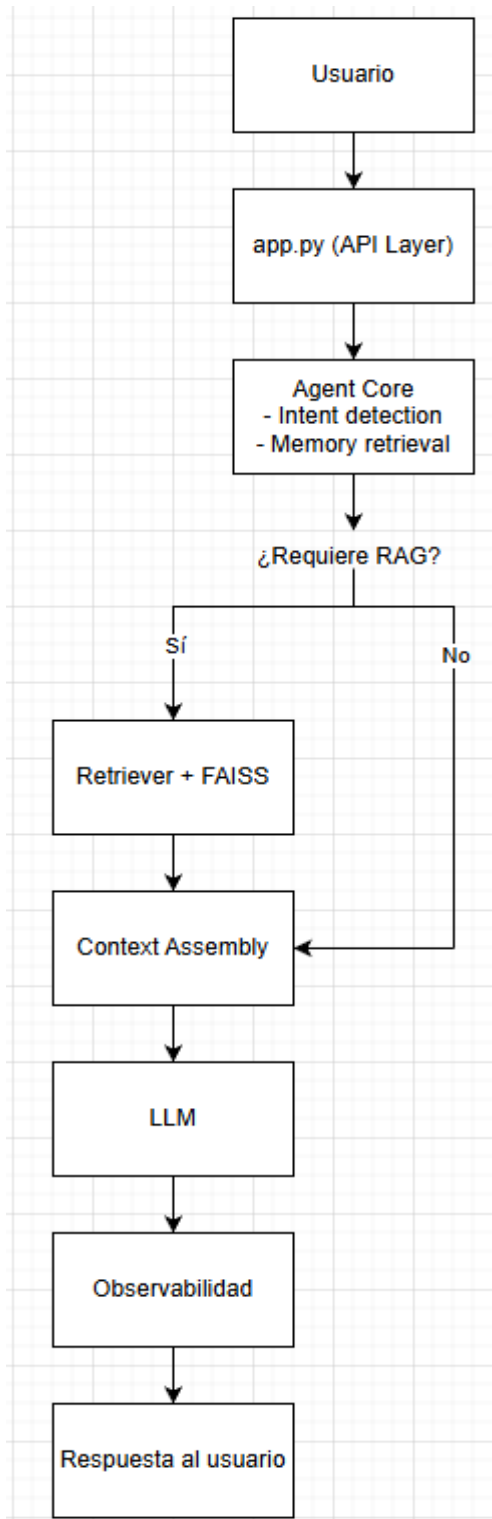
Estrategias de planificación

El agente debe:

- 1. Entender la intención.**
- 2. Clasificar: información → acción → herramienta.**
- 3. Seleccionar si corresponde usar:**
 - ☐ *RAG*
 - ☐ *Memoria*
 - ☐ *Tool*
- 4. Reensamblar contexto.**
- 5. Generar respuesta final con orientación institucional.**

Orquestación

Diagrama textual



Métricas

Latencia

Dividida en:

- Preprocesamiento
- Retrieval
- Reranking
- LLM
Se almacenan timestamps para análisis.

Precisión

Comparación entre:

- Chunks recuperados
- Chunks etiquetados relevantes

Consistencia

Métricas:

- Similitud entre respuestas históricas.
- Detección de variaciones no justificadas.

Logs y análisis

Formato recomendado:

```
timestamp
usuario
consulta_original
embedding_query_id
documentos_recuperados
modelo_utilizado
tiempo_respuesta
resultado
```

Hallazgos:

- *Preguntas ambiguas generan chunks irrelevantes.*
- *La memoria episódica influye en conversaciones largas.*
- *Usuarios intentan pedir información no autorizada — se detectan patrones de seguridad.*

Protocolos de seguridad

- *Sanitizar inputs para evitar ataques como prompt injection.*
- *Clasificar contenido según confidencialidad.*
- *API keys rotadas automáticamente.*
- *Control de longitud del prompt para evitar extracción de información.*

CONCLUSIONES + MEJORAS PROPUESTAS

Conclusiones

IAEva representa una base robusta para construir un asistente institucional, con una arquitectura modular que prioriza:

- *Seguridad*
- *Trazabilidad*
- *Control*
- *Capacidad de expansión*

El sistema demuestra que una organización puede implementar un LLM + RAG sin depender de infraestructuras costosas ni complejas, manteniendo control sobre el conocimiento interno.

Mejoras y propuestas futuras

1. Rendimiento

- *Caching semántico con umbral >0.9.*
- *Precomputación de embeddings.*

2. Memoria

- *Reemplazar JSON con Postgres.*
- *Memoria híbrida basada en embeddings.*

3. Arquitectura

- *Contenedorización con Docker/Kubernetes.*
- *Autoscaling de RAG y LLM.*

4. Seguridad

- *Clasificador de contenido sensible.*

5. Experiencia de usuario

- *Dashboard analítico más completo.*
- *Informe semanal de actividad.*

Reflexiones

Matias Cerda R.: Este proyecto nos mostró que la inteligencia artificial no viene a reemplazar a las personas, sino a potenciar lo que ya hacemos. IAEva no solo responde preguntas; ayuda a liberar tiempo, a reducir la carga operativa y a que cada persona pueda enfocarse en tareas que realmente aportan valor. Entender este equilibrio, donde la tecnología acompaña y no desplaza, es fundamental para adoptarla de manera responsable y consciente.

Javier Muñoz T.: Este proyecto no solo nos enseñó sobre RAG, LLMs y observabilidad; también nos obligó a enfrentar problemas reales de ingeniería: tomar decisiones, lidiar con limitaciones, elegir entre distintos caminos y asumir la responsabilidad de lo que construimos. La inteligencia artificial deja de ser solo un concepto teórico cuando la llevamos a un sistema que interactúa con personas reales, y es ahí donde entendemos el verdadero peso de lo que estamos creando.

Capturas

Configuración

Inicializar Memoria Larga (RAG)

Planificación

Objetivo para planificar

Preparar reporte semanal de consulto

Generar plan

Decidir siguiente paso (demo)

Comandos útiles

guardar: [texto] <guardar texto>

listar notas: <listar notas>

calc: <2+2>

Agente Organizacional – Streamlit + Ollama

Deploy

Chat

Escribe tu mensaje

¿Que se dice en politicas.md?

Enviar

¿A qué hora se entregan los reportes semanales?

Instrucción:

La organización requiere que todos sus miembros sigan un procedimiento específico para realizar la comunicación externa y el onboarding de nuevos empleados. Para asegurarnos de que todos estén al tanto, le proporcionaré los pasos a seguir según nuestros documentos internos.

Comunicación Externa:

Según nuestra política de comunicación (consultado en [nuestro documento de políticas](#)):

- "Toda comunicación externa debe pasar por el área de Comunicaciones."

Esto significa que antes de cualquier acción o mensaje hacia un exterior, debemos asegurarnos de que se haya revisado y aprobado por la unidad correspondiente.

Onboarding:

Para el proceso de incorporación (onboarding) de nuevos empleados, también tenemos unos pasos definidos en [nuestro documento de procedimientos](#):

- "1. Crear cuenta institucional."
- "Asignar roles de acceso."
- "Revisión de políticas el primer día..."

Esto garantiza que todos los empleados nuevos estén informados sobre nuestras políticas y procedimientos desde el primer momento.

Activar Windows
Ve a Configuración para activar Windows.

```

1 [{"timestamp": "2025-11-24T18:30:20.023367", "user_message": "¿A qué hora se entregan los reportes semanales?", "mode": "LLM", "tool": "none", "latency_ms": 90208.7086, "memory_mb": 154.0898}
2 [{"timestamp": "2025-11-24T18:30:25.305570", "user_message": "¿A qué hora se entregan los reportes semanales?", "mode": "RAG", "tool": "vectorstore", "latency_ms": 10086.6326, "memory_mb": 185.3789}
3 [{"timestamp": "2025-11-24T18:32:16.869420", "user_message": "Hablame sobre los reportes semanales", "mode": "RAG", "tool": "vectorstore", "latency_ms": 10086.6326, "memory_mb": 185.3789}
4 [{"timestamp": "2025-11-24T18:33:02.663970", "user_message": "Hablame sobre el procedimiento de onboarding", "mode": "RAG", "tool": "vectorstore", "latency_ms": 10086.6326, "memory_mb": 185.3789}
5 [{"timestamp": "2025-11-24T18:34:03.155695", "user_message": "Hablame sobre el archivo procedimientos.md", "mode": "RAG", "tool": "vectorstore", "latency_ms": 10086.6326, "memory_mb": 185.3789}
6 [{"timestamp": "2025-11-24T18:42:14.550901", "user_message": "Que hora es en Nueva York?", "mode": "LLM", "tool": "none", "latency_ms": 42359.44970000128, "memory_mb": 154.0898}
7 [{"timestamp": "2025-11-24T18:42:22.718235", "user_message": "¿Que hora es en Nueva York?", "mode": "LLM", "tool": "none", "latency_ms": 4736.70630001286, "memory_mb": 154.0898}
8

```



Adjuntos



IAEva-main.zip