# Planning Individual - Desarrollador A (IA y Lógica de Negocio)

## Perfil y Responsabilidades Principales

Rol: AI & Business Logic Engineer

Especialización: Inteligencia Artificial, LLMs, análisis educativo y experiencia de usuario

Áreas de Dominio:

- Integración y optimización de modelos de lenguaje
- Diseño de sistemas de generación de contenido educativo
- Algoritmos de personalización y adaptación
- Análisis de datos educativos
- Experiencia de usuario con IA
- Evaluación y feedback automatizado

## **o** Objetivos Clave del Rol

- 1. Crear el mejor sistema de generación de exámenes adaptativo del mercado
- 2. Optimizar costos de IA en 50% mediante routing inteligente
- 3. Lograr 90% de precisión en evaluaciones automáticas
- 4. Personalización efectiva para cada estudiante
- 5. Tiempo de respuesta <2s para interacciones con IA

## **Planning Detallado por Fases**

FASE 0: Análisis y Diseño de IA (Semana 0)

Día	Tareas	Entregables	Criterios de Éxito
Lunes	• Análisis de LLMs disponibles (GPT-4, Claude, Gemini, Llama) < br>• Comparativa de capacidades y costos < br>• Definir casos de uso por modelo	Matriz comparativa     LLMs de selección	• ROI claro por modelo < br> • Casos de uso mapeados
Martes	Diseñar arquitectura de prompts < br>     Definir templates base < br>     Estrategia de few-shot learning	• Prompt engineering guide • Template library v0.1	Prompts     estructurados < br>     Versionado definido
Miércoles	• Investigar vector databases • Diseñar sistema de embeddings • RAG architecture planning	Comparativa vector  DBs < br>     Arquitectura  RAG	• Solución escalable < br>• < 100ms búsqueda
Jueves	Definir lógica de evaluación Diseñar sistema de scoring aprendizaje	<ul><li>Algoritmos de evaluación &lt; br&gt;</li><li>Rúbricas automatizadas</li></ul>	• Evaluación objetiva < br>• Feedback actionable
Viernes	Workshop con Dev B sobre APIs < br>     Definir flujos de datos < br>     Establecer SLAs de IA	API contracts < br >      Flujos documentados	• Interfaces claras < br>• SLAs realistas

# FASE 1: Integración Base de IA (Semanas 1-2)

## Semana 1: APIs y Modelos Base

Tarea	Descripción Detallada	Dependencias	Riesgos
1.1 Cliente Multi- LLM	• Abstracción para múltiples providers < br> • Gestión de API keys < br> • Retry logic y fallbacks < br> • Rate limiting por provider	API Gateway setup	Límites de rate imprevistos
1.2 Servicio de Autenticación	• JWT implementation < br>• Roles y permisos < br>• Refresh tokens < br>• Session management	Redis disponible	Seguridad de tokens
1.3 Vector  Database Setup	• Pinecone/Weaviate/Qdrant < br > • Índices por tipo de contenido < br > • Estrategia de chunking < br > • Metadata structure	Storage disponible	Costos de escalado
1.4 Prompt Management	Versionado de prompts < br>     A/B testing framework < br>     Variables y templates < br>     Prompt optimization	Base de datos	Gestión de versiones
API Gateway		Complejidad de WS	

## Semana 2: Lógica de Negocio Core

Tarea	Descripción Detallada	Dependencias	Riesgos
2.1 Analizador de Estudiantes	Perfil learning style < br>     Análisis de fortalezas/debilidades < br>     Histórico de rendimiento < br>     Predicciones ML	DB esquemas	Privacidad de datos
2.2 Motor de Personalización	Adaptación de contenido < br>     Dificultad dinámica < br>     Recomendaciones personalizadas < br>     Learning paths	Perfil estudiante	Over- personalización
2.3 Sistema de Embeddings	Generación para documentos Chunking     strategies Semantic search Similarity     scoring	Vector DB	Calidad embeddings
2.4 Cache de Prompts	Respuestas frecuentes < br>     Semantic caching < br>     TTL inteligente < br>     Invalidación selectiva	Redis cluster	Cache coherence
2.5 Analytics Base	• Event tracking < br>• Métricas de uso < br>• Performance metrics < br>• User journey	Analytics DB	Volume de datos

## FASE 2: Generación de Contenido (Semanas 3-4)

### Semana 3: Sistema de Exámenes

Tarea	Descripción Detallada	Skills Necesarios	Métricas de Éxito
3.1 Generador Tipo Test	inteligentes < br>• Balanceo de dificultad < br>• Validación		95% preguntas válidas
3.2 Generador Desarrollo	• Preguntas abiertas < br> • Rúbricas automáticas < br> • Ejemplos de respuesta < br> • Criterios de evaluación	NLP, LLMs	
3.3 Generador Problemas	'		Soluciones correctas
<ul> <li>3.4 Mixer de</li></ul>		Algoritmos	Exámenes balanceados
• Análisis exámenes previos < br>• Extracción de patrones < br>• Aplicación de estilo < br>• Validación formato		ML, Pattern matching	90% similitud estilo
•			<u> </u>

## Semana 4: Evaluación y Feedback

Tarea	Descripción Detallada	Skills Necesarios	Métricas de Éxito
4.1 Corrector Automático	• NLP para respuestas texto < br> • Partial scoring < br> • Detección de conceptos < br> • Explicación de errores	NLP, LLMs 85% accuracy vs humano	
4.2 Generador de Feedback	Feedback personalizado < br >      Sugerencias de mejora < br >      Enlaces a recursos < br >      Motivación adaptativa	Psychology, LLMs	Feedback actionable
4.3 Sistema de Métricas	• Learning analytics < br>• Progress tracking < br>• Predicción de éxito < br>• Early warning system	ML, Statistics	Predicciones 80% accurate
4.4 Recomendador de Estudio	Gaps identification < br>     Material prioritization < br>     Schedule optimization < br>     Adaptive planning	ML algorithms	Mejora 20% resultados

# FASE 3: Optimización de IA (Semanas 5-6)

## Semana 5: Routing y Optimización

Tarea	Descripción Detallada	Herramientas	Objetivo
5.1 Router Inteligente	Clasificación de queries < br > • Modelo selection     logic < br > • Cost/quality balance < br > • Fallback strategies	RouteLLM, ML 50% reducción costo	
5.2 Fine-tuning Models	Dataset preparation < br>     Training pipelines < br>     Evaluation metrics < br>     Model versioning	OpenAl API, +15% accuracy	
5.3 Prompt Optimization	Automatic prompt tuning < br > • A/B testing     results < br > • Template refinement < br > • Context     optimization	DSPy, LangChain	-30% tokens used
<ul> <li>5.4 Batch</li> <li>Async generation &lt; br&gt;         <ul> <li>Priority queues &lt; br&gt;             <li>Bulk</li> <li>operations &lt; br&gt;             <li>Progress tracking</li> </li></li></ul> </li> </ul>		Celery, Redis	10x throughput
• Error recovery < br > • Content filtering < br > • Safety  Handling checks < br > • Quality assurance		LLM guards	<1% bad outputs

## Semana 6: Caché y Performance

Tarea	Descripción Detallada	Herramientas	Objetivo
6.1 Semantic Cache	• Embedding-based cache < br> • Similarity threshold < br> • Cache warming < br> • Hit rate optimization	Pinecone, Redis	70% cache hit
6.2 Response Streaming	• SSE implementation < br> • Chunked responses < br> • Progressive rendering < br> • Perceived performance	FastAPI, SSE	<500ms first byte
6.3 Precomputed Responses	• Common questions DB br>• Template responses • Quick replies br>• Instant feedback	PostgreSQL	<50ms common queries
6.4 Model  Quantization	• Local model optimization < br>• Quantization testing < br>• Performance benchmarks < br>• Accuracy trade-offs	GGUF, llama.cpp	2x speed, -5% acc

# **FASE 4: Features Avanzadas (Semanas 7-8)**

## Semana 7: Agentes y Automación

Tarea Descripción Detallada Comp		Complejidad	Prioridad
7.1 Agente de	Multi-step reasoning < br>     Tool usage (calculator, search)	Alta	P0
Estudio		Alla	PU
7.2 Tutor	Context-aware chat Socratic method Adaptive	Alta	PO
Conversacional	responses < br>• Emotion detection	Alla	PU
7.3 Auto-curriculum	• Learning path generation < br>• Prerequisite mapping < br>•	Media	D1
7.5 Auto-curriculum	Pacing optimization • Milestone setting	Media	P1
7.4 Collaborative	Peer matching < br > • Group study sessions < br > • Shared	N 41: -	D2
Learning	resources < br>• Competition elements	Media	P2
4	'	•	<b>•</b>

### Semana 8: Análisis Avanzado

Tarea	Tarea Descripción Detallada		Prioridad
8.1 Predictive	• Success prediction < br> • Risk identification < br> • Intervention	Alta	D0
Analytics	triggers < br>• Outcome modeling	Alla	P0
8.2 Knowledge	Concept mapping < br > • Relationship extraction < br > •	Alta	P1
Graphs	Prerequisite chains < br>• Gap analysis	Alla	PI
8.3 Adaptive	<b>8.3 Adaptive</b> • CAT implementation < br>• IRT models < br>• Real-time		D1
Testing	adaptation < br>• Precision stopping	Alta	P1
<b>8.4 Learning</b> • Deep analytics • Pattern recognition •		Media	P2
Insights Recommendations engine < br> • Parent dashboards		ivieuia	F Z
4			•

## FASE 5: Innovación y Escala (Semanas 9-10)

### **Semana 9: Características Experimentales**

Tarea Descripción Detallada Stack Téc		Stack Técnico	Innovación
9.1 Voice	Speech-to-text Oral examinations	A SI SI	
Interaction	Pronunciation feedback • Audio lessons	Whisper API	Accessibility
9.2 AR/VR	• 3D visualizations • Immersive learning • Spatial	Three.js,	
Integration	memory < br>• Gamification	WebXR	
9.3 Multi-modal	Image understanding < br >      Diagram generation < br >	GPT-4V,	
Learning	Video analysis < br>• Mixed media	DALL-E	Comprehension
9.4 Peer Learning	Student matching < br>     Collaborative AI < br>     Group	Graph     Retention	
AI	dynamics < br> • Social learning		
■	•	•	•

### Semana 10: Preparación para Escala

Tarea	Descripción Detallada Stack Técnico Impa		Impacto
10.1 A/B Testing	Feature flags Experiment design	Laura de Davido	Data division
Platform	Statistical significance < br> • Rollout strategies	LaunchDarkly	Data-driven
10.2 MI Dinalina	Feature engineering < br>     Model training < br>	MLflow,	Continuous
10.2 ML Pipeline	Deployment automation < br> • Monitoring	Kubeflow	improvement
10.3 API	OpenAPI specs < br > • Interactive docs < br > • Code	Swagger,	Developer
Documentation	examples < br>• SDKs generation	Postman	adoption
10.4 Knowledge	Technical documentation < br>     Architecture		
10.4 Knowledge	decisions < br>• Best practices guide < br>• Training	Confluence	Team scalability
Transfer	sessions		
◀			<u> </u>

## **III** KPIs y Métricas de Éxito

#### Métricas de IA

• Accuracy de evaluación: >85% vs evaluación humana

• Costo por usuario: <\$0.50/mes en LLM

• Latencia generación: <2s para respuestas

• Cache hit rate: >70% consultas similares

### Métricas de Producto

• Satisfacción usuario: NPS > 50

• **Engagement**: 5+ sesiones/semana por usuario

• Mejora académica: +20% en calificaciones

• **Retención**: >80% usuarios activos mensual

#### Métricas de Calidad

• Calidad de exámenes: 95% válidos sin revisión

Personalización efectiva: 90% usuarios reportan mejora

• Tiempo ahorro profesor: 10h/semana

Cobertura curricular: 100% temas

## K Stack Tecnológico Principal

#### AI/ML Stack

• **LLMs**: OpenAl, Anthropic, Google, Open models

Embeddings: OpenAl Ada, Sentence Transformers

Vector DB: Pinecone/Weaviate/Qdrant

ML Framework: PyTorch, Transformers

#### **Backend Framework**

• API: FastAPI + Pydantic

• **Async**: asyncio, aiohttp

Task Queue: Celery + Redis

WebSockets: Socket.io

#### **AI Tools**

Orchestration: LangChain/LlamaIndex

• **Optimization**: RouteLLM, vLLM

• **Evaluation**: Ragas, LangSmith

Fine-tuning: OpenAl API, Axolotl

### **Data & Analytics**

• Analytics: Mixpanel/Amplitude

Experiments: Statsig/LaunchDarkly

Monitoring: Weights & Biases

• Data Pipeline: Apache Beam

## 🚀 Estrategias de Optimización

#### Reducción de Costos LLM

#### 1. Clasificación de Queries

- Simple → Modelo pequeño/caché
- Complejo → GPT-4/Claude
- Creativo → Modelo especializado

#### 2. Caché Inteligente

- Semantic similarity matching
- Template-based responses
- Precomputed common queries

### 3. Batch Processing

- Agrupar requests similares
- Off-peak processing
- Bulk generation

### Mejora de Calidad

#### 1. Fine-tuning Continuo

- Feedback loop de usuarios
- Dataset curation
- Model evaluation

#### 2. Prompt Engineering

- A/B testing prompts
- Few-shot examples
- Chain-of-thought

#### 3. Ensemble Methods

- Multiple model voting
- Confidence scoring
- Quality filtering

## Algoritmos Clave a Implementar

#### Sistema de Evaluación

python

```
# Pseudocódigo conceptual

class ExamEvaluator:

def evaluate_answer(self, question, student_answer, rubric):

# 1. Extraer conceptos clave

# 2. Comparar con rúbrica

# 3. Partial credit calculation

# 4. Generar feedback específico

pass
```

### Personalización Adaptativa

```
python

# Pseudocódigo conceptual

class AdaptiveLearning:

def next_content(self, student_profile, performance_history):

# 1. Calcular nivel actual

# 2. Identificar gaps

# 3. Seleccionar siguiente contenido

# 4. Ajustar dificultad

pass
```

#### **Router de Modelos**

```
python

# Pseudocódigo conceptual

class ModelRouter:

def route_query(self, query, context, constraints):

# 1. Clasificar tipo de query

# 2. Estimar complejidad

# 3. Check caché

# 4. Seleccionar modelo óptimo

pass
```

# Riesgos y Mitigaciones Específicos

Riesgo	Probabilidad	Impacto	Mitigación
Costos LLM explotan	Alta	Alto	Límites duros, alertas, caché agresivo
Hallucinations en exámenes	Media	Alto	Validation layers, human review option
Sesgo en evaluaciones	Media	Alto	Bias testing, diverse training data
Latencia inaceptable	Baja	Alto	Streaming, perceived performance tricks
Cambios en APIs LLM	Media	Medio	Abstraction layer, multiple providers
4		•	•

### Plan de Formación Continua

### **Cursos y Certificaciones**

- 1. Fast.ai Practical Deep Learning
- 2. Stanford CS224N (NLP)
- 3. Anthropic's Constitutional Al
- 4. Google's Machine Learning Crash Course

### **Investigación y Papers**

- Attention mechanisms y transformers
- Few-shot y zero-shot learning
- Retrieval Augmented Generation
- Educational technology research

### **Comunidad y Networking**

- Participar en Hugging Face community
- Contribuir a open source LLM projects
- Asistir a conferencias AI/EdTech
- Blog técnico sobre learnings

## **o** Hitos Clave por Fase

#### **Fase 1: Foundation**

- ✓ Primera generación de examen exitosa
- ✓ API respondiendo <2s</li>
- ✓ 3 tipos de LLM integrados

#### **Fase 2: Core Features**

- ✓ 95% accuracy en tipo test
- ✓ Feedback personalizado funcionando
- ✓ 1000 exámenes generados

### **Fase 3: Optimization**

- ✓ 50% reducción en costos
- ✓ 70% cache hit rate
- ✓ <500ms perceived latency</li>

#### Fase 4: Advanced

- ✓ Agente tutor conversacional
- ✓ Predicciones 80% accurate
- ✓ Multi-modal support

#### Fase 5: Scale

- ✓ A/B testing platform live
- ✓ ML pipeline automatizado
- ✓ 100% documentación

## 🔽 Definition of Done para IA

Para cada feature de IA:

- Prompt optimizado y versionado
- Evaluación contra ground truth
- Métricas de calidad definidas
- Fallbacks implementados
- Costos calculados y optimizados
- Tests de edge cases
- Documentación de limitaciones
- Monitoring en producción

### Puntos de Sincronización con Dev B

#### Semanal

- API contracts review
- Performance bottlenecks
- Infrastructure needs
- Cost optimization

#### **Quincenal**

- Architecture decisions
- Integration testing
- Deployment planning
- Retrospectiva técnica

#### Mensual

Roadmap alignment

- Technical debt review
- Innovation planning
- Knowledge sharing