

# Especificación del Módulo de Fine-Tuning

---

## 1. Introducción

**1.1 Objetivo del Módulo** El módulo de *fine-tuning* permite a usuarios con conocimientos básicos realizar Supervised Fine-Tuning (SFT) y el uso de técnicas de PEFT (Parameter Efficient Fine Tuning, adaptadores) en modelos de lenguaje pequeños (menos de 15 mil millones de parámetros). Su propósito es simplificar el proceso, reducir la complejidad técnica y optimizar el uso de recursos computacionales, enfocándose en escenarios de bajo costo y alta accesibilidad. Está diseñado como una base extensible para soportar técnicas avanzadas y modelos más grandes en futuras iteraciones.

### 1.2 Alcance Inicial

- **Técnicas:** SFT con soporte opcional para adaptadores LoRA y QLoRA (LoRA cuantizado).
- **Modelos:** Modelos de lenguaje con menos de 15B parámetros.
- **Infraestructura:** Ejecución en infraestructura propia del usuario.
- **Interfaz:** Configuración sencilla de datasets y parámetros básicos para usuarios no expertos.
- **Subproductos:** Métricas, logs, eventos y checkpoints generados durante el proceso.

### 1.3 Limitaciones

- No soporta modelos grandes (>15B parámetros).
- Sin integración con proveedores externos de *fine-tuning* o despliegue remoto.
- Excluye técnicas como *Reinforcement Learning* (RL) o destilación de modelos.
- Requiere infraestructura local con recursos suficientes (CPU/GPU).

Estas limitaciones serán levantadas en versiones futuras.

---

## 2. Requisitos

### 2.1 Requisitos Funcionales

#### 2.1.1 Configuración del Fine-Tuning

- **Datasets:**
  - Entrenamiento: Obligatorio (CSV/JSON, máximo algunos miles de registros, 1GB).

- Evaluación: Opcional (CSV/JSON, unos pocos miles de registros, 1GB).
- Formato:
  - \* SFT: { "text": str, "label": str }.
  - \* LoRA/QLORA: Igual que SFT.
- **Hiperparámetros:** Se soportará en un principio los parámetros básicos:
 

Parámetro	Tipo	Rango/Default	Descripción
epochs	Entero	1-10 / 3	Número de épocas de entrenamiento.
learning_rate	Float	1e-6 - 1e-3 / 5e-5	Tasa de aprendizaje inicial.
batch_size	Entero	1-256 / 32 o "auto"	Tamaño del batch (auto ajusta según GPU)
lora_rank	Entero	4-64 / 16	Rango de LoRA (si aplica).
lora_alpha	Float	8-64 / 32	Factor de escalado de LoRA (si aplica).
lora_dropout	Float	0-0.5 / 0.1	Probabilidad de dropout de LoRA.
quantization_bits	Entero	4 o 8 / 8	Bits para QLORA (si aplica).
- **Configuración Avanzada:** Archivo JSON opcional con parámetros adicionales (e.g., {"warmup\_steps": 100}).

**2.1.2 Gestión de Datasets** Se reconoce la necesidad de permitir el registro de datasets, estos datasets pueden ser públicos o privados. - **Acciones:** Registro, carga, eliminación. - **Tipos:** Públicos (acceso general) y privados (por organización). - **Validación:** Métrica de completitud (% de campos no nulos). - **Almacenamiento:** Límite de 1GB por dataset, con cuotas por usuario (TBD). A futuro se debe tener la capacidad de trabajar con datasets, generación con datos ingresados, de traza, datos sintéticos, edición y validación.

**2.1.3 Modelos y Adaptadores** Se debe llevar registro de todos los modelos y adaptadores, públicos y privados.

- **Modelos:**
  - Registro de modelos base y afinados.
  - Metadatos: { "id": str, "name": str, "type": "base|finetuned", "origin": "local|remote", "created\_at": datetime }.
- **Adaptadores:**
  - Pesos entrenados para tareas/dominios específicos.
  - Merge en memoria o combinación con modelos base.
- **Privacidad:** Modelos y adaptadores restringidos por organización.

#### 2.1.4 Despliegue

- **Local:** Despliegue en motor de inferencia propio (e.g., PyTorch, vLLM, slang, ollama).
- **Remoto:** No soportado en esta fase.

#### 2.1.5 Trabajos de Fine-Tuning

- **Definición:** Cada tarea es un *job* asíncrono.

- **Subproductos:**
  - **Métricas:** Pérdida de entrenamiento (`train_loss`), pérdida de validación (`valid_loss`), precisión (si aplica).
  - **Tiempos:** Duración total, tiempo por época.
  - **Recursos:** Uso de CPU/GPU, memoria consumida.
  - **Logs:** Registro detallado de ejecución.
  - **Eventos:** Notificaciones (inicio, fin, errores).
  - **Checkpoints:** Guardado cada 500 pasos (configurable).

## 2.2 Requisitos No Funcionales

- **Rendimiento:** Tiempo de respuesta del API  $< 2$  segundos.
- **Seguridad:** Autenticación por organización; datasets y modelos privados no accesibles fuera de ella.
- **Compatibilidad:** Python 3.8+, bibliotecas Hugging Face (`transformers`, `datasets`, `peft`).
- **Escalabilidad:** Soporte de varios trabajos concurrentes por usuario.

## 3. Arquitectura

### 3.1 Módulos de Alto Nivel

1. **FineTuner:** Orquesta trabajos de *fine-tuning*.
2. **Datasets:** Administra carga, filtrado y evaluación de datos.
3. **Models:** Gestiona modelos y adaptadores.
4. **Deployment:** Controla despliegue local.
5. **Evals:** Gestión de evaluaciones.

El usuario entrenará mediante fine tuning modelos o adaptadores. Para ello someterá tareas de fine tuning de acuerdo a los parámetros seleccionados. Para ello deberá utilizar datasets de entrenamiento y validación. Luego podrá desplegar estos modelos con o sin adaptadores. Sobre estos modelos desplegados podrá realizar inferencias o evaluaciones.

Módulo	Función
<b>FineTuner</b>	Realiza el fine tuning.
<b>Dataset</b>	Gestión de datasets.
<b>Model</b>	Gestión de modelos.
<b>Adapter</b>	Gestión de adaptadores de un modelo.
<b>Checkpoint</b>	Gestión de checkpoints de un modelo.
<b>Deployment</b>	Gestión de deployments.
<b>Evals</b>	Gestión de evaluaciones

**3.2 Módulos Internos** Se requiere soporte de los siguientes módulos para el desarrollo del módulo.

Módulo	Función
<b>Files</b>	Manejo de archivos (datasets, logs, etc.).
<b>Jobs</b>	Gestión de trabajos asíncronos.
<b>Logs</b>	Registro de actividades.
<b>Metrics</b>	Cálculo de métricas (e.g., pérdida, calidad).
<b>Monitoring</b>	Supervisión de recursos y progreso.
<b>Events</b>	Generación de notificaciones.

**3.3 Diagrama de Relación** Se reconocen los siguientes niveles a desarrollar.

[API Externa] [SDK]

[FineTuner] [Datasets] [Models] [Adapters] [Evals] [Checkpoints] [Deployment]

[Jobs] [Files, Logs, Metrics, Monitoring, Events]

### 3.4 Diseño Agnóstico

- Interfaces abstractas para soportar infraestructura propia o externa.
- Configuración basada en JSON para extensibilidad.

## 4. Interfaces

**4.1 API Externa (OpenAPI)** Se recomienda en este punto analizar la API de Open API para fine tuning. - **/fine\_tuning/jobs:** - POST /create: { "base\_model": str, "dataset": str, "hyperparameters": dict, "output\_model": str } → { "job\_id": str }. - GET /list: [ { "job\_id": str, "status": "pending|running|succeeded|failed|cancelled", "created\_at": datetime } ]. - **/fine\_tuning/jobs/{job\_id}:** - GET /: { "job\_id": str, "status": str, "progress": float, "metrics": dict, "result": dict }. - POST /cancel: Cancela el trabajo. - GET /checkpoints: [ { "id": str, "step": int, "metrics": dict } ]. - GET /events: [ { "id": str, "level": "info|error", "message": str } ]. - **/deployment/model:** - POST /deploy: { "model\_id": str, "adapters": [str] } → { "deployment\_id": str }. - GET /{model\_id}: { "deployment\_id": str, "status": str }. - POST /undeploy: Termina el despliegue.

### 4.2 SDK (Línea de Comandos)

Comando	Descripción
list datasets	Lista datasets registrados.

Comando	Descripción
<code>create dataset</code> <code>&lt;name&gt; &lt;path&gt;</code>	Registra un dataset desde archivo.
<code>delete dataset</code> <code>&lt;name&gt;</code>	Elimina un dataset.
<code>list models</code>	Lista modelos registrados.
<code>create model &lt;name&gt;</code> <code>&lt;path&gt;</code>	Registra un modelo desde archivos locales.
<code>deploy &lt;model-name&gt;</code> <code>[-adapters]</code>	Despliega un modelo con adaptadores opcionales.
<code>create</code> <code>fine-tune-job</code> <code>[options]</code>	Crea un trabajo con opciones: <code>-basemodel</code> , <code>-dataset</code> , <code>-epochs</code> , <code>-lora-rank</code>
<code>get fine-tuning-job</code> <code>&lt;job-name&gt;</code>	Obtiene estado y métricas del trabajo.

Ejemplo:

```
create fine-tune-job -basemodel "bert-base-uncased" -dataset "qa_data.csv" -output-model "be
```

## 5. Evolución Futura

### 5.1 Roadmap

1. **Fase 1:** SFT + LoRA/QLORA en infraestructura propia.
2. **Fase 2:** RL Fine-Tuning y modelos >15B parámetros.
3. **Fase 3:** Soporte para proveedores externos y jerarquías de modelos/adaptadores.
4. **Fase 4:** Destilación y generación de datos sintéticos.

### 5.2 Consideraciones

- **Extensibilidad:** Interfaces modulares para nuevos métodos de *fine-tuning*.
- **Monetización:** Posibilidad de ofrecer como servicio SaaS.
- **Organización:** Árbol jerárquico de modelos/adaptadores por tarea/dominio.

## 6. Entidades

Entidad	Atributos
FineTuningJob	id: str, organization_id: str, created_at: datetime, status: str, base_model: str, dataset: str, hyperparameters: dict, result_files: [str]
Checkpoint	id: str, job_id: str, step_number: int, metrics: { "train_loss": float, "valid_loss": float }
Event	id: str, job_id: str, level: "info error", message: str, created_at: datetime
File	id: str, path: str, type: "dataset model log"
Model	id: str, name: str, type: "base fine-tuned", origin: "local remote", created_at: datetime
Adapter	id: str, model_id: str, dataset: str, created_at: datetime, parameters: dict

## 7. Ejemplo de Uso

**Escenario:** Afinar un modelo para QA. 1. Registrar dataset: `create dataset qa_data qa_data.csv`. 2. Crear trabajo: `create fine-tune-job -basemodel "bert-base-uncased" -dataset "qa_data" -output-model "bert-qa" -epochs 3`. 3. Consultar estado: `get fine-tuning-job <job-id>`. 4. Desplegar: `deploy bert-qa`.

## 8. Diagramas

Diagrama componentes:

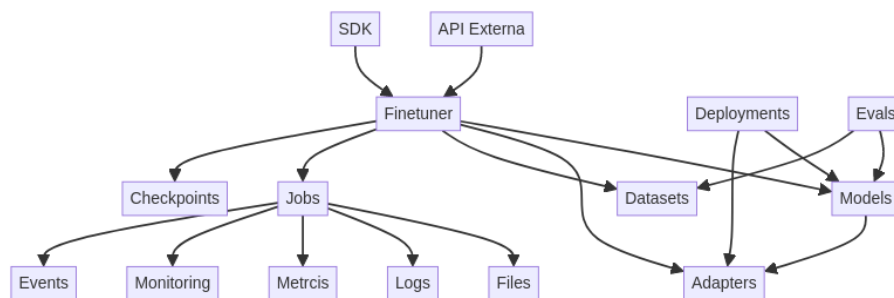


Diagrama módulos:

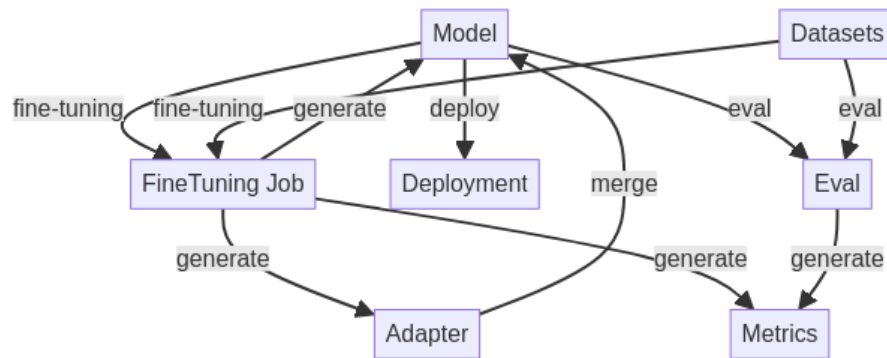


Diagrama entidades:

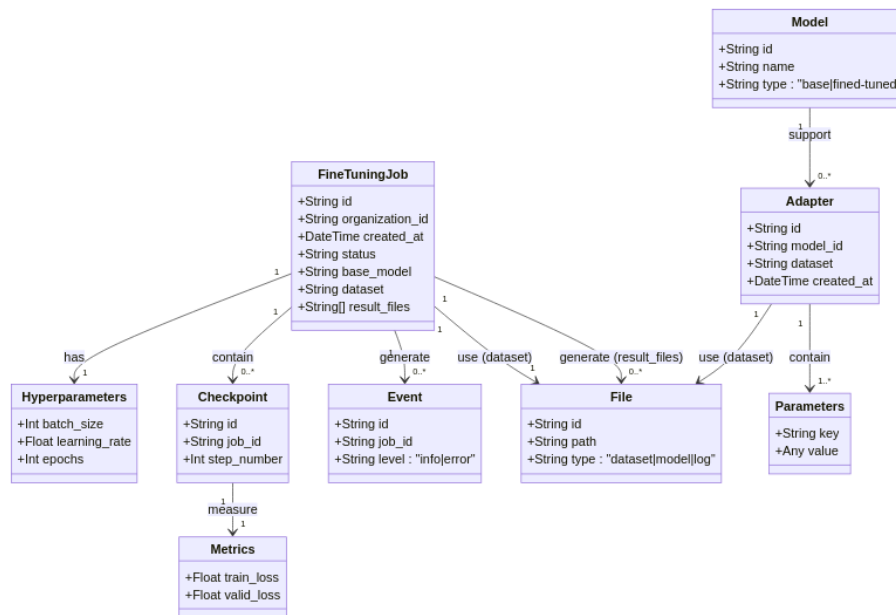


Diagrama de API de TEST:

