

# Document Technique : Avikam1 LLM

*Un modèle de langage efficace pour Evilafo AI*

## Résumé

*Avikam1* est un modèle de langage (LLM) efficace et optimisé, conçu spécifiquement pour les applications bilingues français-anglais, notamment le chatbot *EvilaFo AI*. Basé sur une architecture *Transformer* améliorée, ce modèle combine **LoRA** (*Low-Rank Adaptation*) et **FlashAttention 2** pour réduire significativement les coûts computationnels tout en maintenant des performances compétitives.

Nos contributions principales incluent :

- Une **méthode de fine-tuning économe en ressources**, divisant par trois l'utilisation de mémoire GPU grâce à des adaptateurs légers (réduction de 68% en VRAM).
- Une **optimisation pour le français**, via un prétraitement spécifique et un fine-tuning sur des corpus bilingues.
- Un **déploiement industrialisable**, avec une API REST à faible latence (< 50 ms) et une intégration transparente avec *MLflow* pour le suivi des expériences.

Les benchmarks démontrent qu'\*Avikam1-7B\* atteint une **perplexité de 12.3** sur des textes français (contre 14.1 pour LLaMA2-7B), avec un coût d'entraînement réduit de 57% (1 200vs2800vs2800). Ces avancées en font une solution idéale pour les applications nécessitant des LLMs performants mais économiques, comme les assistants conversationnels ou les outils de génération de texte.

# Sommaire

## 1. Introduction

- 1.1 Contexte et Motivations
- 1.2 Objectifs Scientifiques
- 1.3 Contributions Majeures

## 2. Architecture du Modèle

- 2.1 Transformers Adaptés
- 2.2 Mécanisme d'Attention Optimisé
- 2.3 Intégration de LoRA
- 2.4 Configuration des Paramètres

## 3. Entraînement et Optimisation

- 3.1 Pipeline de Données
- 3.2 Stratégies de Fine-Tuning
- 3.3 Fonction de Coût et Régularisation
- 3.4 Optimisation Mémoire

## 4. Résultats Expérimentaux

- 4.1 Métriques de Performance
- 4.2 Benchmark Comparatif
- 4.3 Analyse des Coûts Computationnels

## 5. Déploiement et Industrialisation

- 5.1 API REST et Latence
- 5.2 Monitoring avec MLflow
- 5.3 Intégration dans EvilaFo AI

## 6. Conclusion et Perspectives

6.1 Bilan des Résultats

6.2 Limitations Actuelles

6.3 Voies d'Amélioration Futures

# Introduction

## Contexte et Motivations

L'essor récent des *Large Language Models* (LLMs) a révolutionné le domaine du traitement automatique du langage (TAL), avec des modèles comme GPT-4, LLaMA-2 et Mistral démontrant des capacités impressionnantes. Cependant, leur déploiement dans des applications spécifiques, telles que les chatbots, se heurte à des défis majeurs :

- **Coût computationnel prohibitif** (e.g., 175B paramètres pour GPT-3)
- **Latence élevée** en environnement de production
- **Adaptation limitée** aux langues autres que l'anglais

Dans ce contexte, **Avikam1** a été conçu pour combler ces lacunes, en proposant un modèle **efficace, spécialisé pour le français et l'anglais**, et optimisé pour des applications temps réel comme EvilaFo AI.

## Objectifs Scientifiques

Nos travaux visent à répondre à trois enjeux clés :

### 1. Efficacité :

Réduire la complexité computationnelle via une architecture hybride **LoRA (Low-Rank Adaptation) + FlashAttention** :

$$\text{VRAM}_{\text{usage}} \propto \underbrace{4 \cdot d^2}_{\text{Paramètres originaux}} + \underbrace{2 \cdot r \cdot d}_{\text{Adaptateurs LoRA}}$$

où ( $d$ ) est la dimension des embeddings

et ( $r \ll d$ ) le rang des adaptateurs (typiquement ( $r = 64$ )).

### 2. Performance Linguistique :

- Optimisation pour le **français** via un pré-traitement spécifique :

```
tokenizer.add_special_tokens({"additional_special_tokens": ["[FR]"]})
```

- Finetuning sur des corpus bilingues.

### 3. Intégration Industrielle :

- Déploiement via une **API REST légère** (< 50 ms de latence) :

```
tokenizer.add_special_tokens({"additional_special_tokens": ["[FR]"]})
```

```
curl -X POST "http://api.evilafo.xyz/generate" -d '{"prompt":"Bonjour", "max_tokens":50}'
```

## Contributions Majeures

Ce document présente les avancées suivantes :

- **Une architecture novatrice** combinant :
  - **LoRA évolutif** pour un fine-tuning efficace (jusqu'à **×3 moins de VRAM**).
  - **FlashAttention 2** pour une accélération de l'inférence :

$$\text{FLOPs}_{\text{attention}} \propto N \cdot d \cdot r \quad (\text{au lieu de } N^2 \cdot d)$$

où  $NN$  est la longueur de la séquence.

- **Un pipeline de fine-tuning optimisé** :
  - Augmentation de données via **back-translation** pour les paires FR/EN.
  - Sélection dynamique des couches à adapter via **analyse du gradient**.
- **Benchmarks complets** contre les baselines (LLaMA-2, Mistral) :

Modèle	Perplexité (FR)	Coût Entraînement (\$)
Avikam1-7B	12.3	1,200
LLaMA2-7B	14.1	2,800

## 2. Architecture du Modèle

### 2.1 Transformers Adaptés

L'architecture de base d'Avikam1 repose sur une variante optimisée du modèle Transformer, spécialement conçue pour équilibrer performance et efficacité. Les principales innovations incluent :

- **Couches Feed-Forward Modifiées :**

python

Copy

Download

```
class AvikamFFN(nn.Module):
    def __init__(self, config):
        super().__init__()
        self.dense1 = nn.Linear(config.hidden_size, config.intermediate_size)
        self.dense2 = nn.Linear(config.intermediate_size, config.hidden_size)
        self.activation = nn.GELU(approximate='tanh') # GELU optimisé
```

Avec un facteur d'expansion réduit à 2.5x (vs 4x classique) pour diminuer les paramètres.

- **Normalisation RMS** (Root Mean Square) :

$$\text{RMSNorm}(x) = \frac{x}{\sqrt{\text{mean}(x^2) + \epsilon}} \odot g$$

Où ( $g$ ) sont les paramètres appris, réduisant le temps de calcul de 15% par rapport au LayerNorm standard.

### 2.2 Mécanisme d'Attention Optimisé

L'implémentation utilise **FlashAttention v2** pour une efficacité maximale :

1. **Attention Hybride :**

- Combinaison d'attention globale (pour les tokens clés) et locale (fenêtre de 256 tokens)

Réduction de la complexité mémoire de  $(O(N^2))$  à  $(O(N^{1.5}))$

## 2. Projections des Clés/Valeurs Groupées :

```
class GroupedQKV(nn.Module):  
    def __init__(self, config):  
        super().__init__()  
        self.qkv = nn.Linear(config.hidden_size, 3*config.hidden_size)  
        self.num_groups = 4 # Projections parallèles
```

### 2.3 Intégration de LoRA

L'adaptation par Low-Rank Adaptation est implémentée avec :

- **Initialisation Orthogonale :**

$$A \sim \mathcal{N}(0, \sigma^2), \quad B = 0$$

Avec  $(\sigma = 1/\sqrt{r})$  pour stabiliser l'apprentissage.

- **Modules Cibles :**

Module	Ratio Paramètres
q_proj	38%
v_proj	42%
o_proj	20%

### 2.4 Configuration des Paramètres

Le tableau récapitule l'architecture de base :

Paramètre	Valeur	Description
hidden_size	4096	Dimension des embeddings
num_layers	32	Couches Transformer

Paramètre	Valeur	Description
num_heads	32	Têtes d'attention
lora_rank	64	Rang des adaptateurs
vocab_size	32000	Taille du tokenizer



## 3. Entraînement et Optimisation

### 3.1 Pipeline de Données

#### Prétraitement :

##### 1. Nettoyage :

- Filtrage des caractères non-UTF8
- Normalisation Unicode (NFKC)

```
def clean_text(text):  
    text = unicodedata.normalize('NFKC', text)  
    return re.sub(r'^\w\s', '', text)
```

##### 2. Tokenization Adaptative :

- Utilisation d'un tokenizer SentencePiece entraîné sur un corpus français/anglais
- Taille de vocabulaire : 32,000 tokens

```
tokenizer = AutoTokenizer.from_pretrained(  
    "evilafo/avikam1-tokenizer",  
    extra_special_tokens=["[FR]", "[EN]"]  
)
```

#### Stratification :

- Répartition 80/10/10 (train/val/test)
- Échantillonnage stratifié par domaine (technique, conversationnel, etc.)

### 3.2 Stratégies de Fine-Tuning

#### Approche en Deux Phases :

##### 1. Phase 1 (3 époques) :

- Taux d'apprentissage : 5e-5
- Batch size : 64
- Gel des couches basses (1-16)

##### 2. Phase 2 (2 époques) :

- Taux d'apprentissage : 1e-5

- Batch size : 32
- Dégel progressif des couches

### Optimiseur :

- AdamW avec rectification de poids (AdamWR)

```
optim:
  name: "adamw"
  lr: 5e-5
  weight_decay: 0.01
  betas: [0.9, 0.999]
  eps: 1e-8
```

### 3.3 Fonction de Coût et Régularisation

#### Loss Composite :

$$\mathcal{L} = \mathcal{L}_{\mathcal{CE}} + \lambda_1 \mathcal{L}_{\mathcal{KL}} + \lambda_2 \|\theta\|_2$$

- -  $\mathcal{L}_{\mathcal{CE}}$  : Cross-entropy classique
- -  $\mathcal{L}_{\mathcal{KL}}$  : Divergence KL pour le lissage des logits
- -  $\lambda_1 = \lambda_2 = 0.0$

#### Dropout Adaptatif :

- Taux initial : 0.1
- Ajustement dynamique basé sur la perplexité de validation

### 3.4 Optimisation Mémoire

#### Techniques Clés :

1. **Gradient Checkpointing :**
  - Activation sélective sur les couches 8, 16, 24, 32
  - Réduction mémoire : ~65%
2. **Quantification 8-bit :**

```

model = AutoModelForCausalLM.from_pretrained(
    "evilafo/avikam1-7b",
    load_in_8bit=True,
    device_map="auto"
)

```

### 3. Packaging des Gradients :

- Regroupement des gradients par blocs de 512MB
- Communication asynchrone NCCL

### Benchmark Mémoire :

Technique	VRAM (GB)
Baseline	24.6
+ LoRA	18.2
+ 8-bit	10.5
+ Checkpointing	7.8

### 3.5 Planification de l'Apprentissage

#### Cyclage des Taux :

- Schedule triangulaire :

$$lr = \text{base\_lr} * (1 - t/T) + \text{min\_lr} * (t/T)$$

Où T = nombre total d'itérations

#### Warmup Progressif :

- 10% des steps pour la phase 1
- 5% pour la phase 2

## 4. Résultats Expérimentaux

### 4.1 Métriques de Performance

#### Benchmark Multilingue :

Nous évaluons Avikam1-7B sur 5 tâches standardisées :

Tâche	FR (Accuracy)	EN (Accuracy)	Δ (FR-EN)
Compréhension (FQuAD)	78.2%	75.1%	+3.1%
Génération (BLEU)	32.4	29.7	+2.7
Raisonnement (LogiQA)	63.5%	61.2%	+2.3%

#### Formule de Score Composite :

Score = 0.4 × Accuracy + 0.3 × BLEU + 0.3 × Perplexité<sub>norm</sub>

### 4.2 Benchmark Comparatif

#### Comparaison Modèles 7B (sur jeu de test FR) :

Modèle	Perplexité	VRAM (GB)	Latence (ms)	Coût Entraînement (\$)
Avikam1-7B	12.3	7.8	42	1,200
LLaMA2-7B	14.1	24.6	58	2,800
Mistral-7B	13.2	22.1	47	2,500

#### Avantages Clés :

- 3.2× moins de mémoire que LLaMA2
- 27% plus rapide en inférence
- Coût réduit de 57%

### 4.3 Analyse des Coûts Computationnels

#### Breakdown du Coût (pour 1M tokens) :

Composant	Coût (GPUh)	% Total
Forward Pass	1.2	38%
Backward Pass	2.1	52%
Overhead LoRA	0.3	10%

Équation d'Efficacité :

$$\text{Efficacité} = \frac{\text{Tokens Traités}}{\text{GPUh}} \times \frac{1}{\text{Perplexité}}$$

Avikam1: **142** vs LLaMA2: **89**

4.4 Analyse Qualitative

Exemple de Sortie :

Prompt: "Explique le théorème de Pythagore en termes simples"

Réponse Avikam1:  
"Dans un triangle rectangle, le carré de l'hypoténuse (côté opposé à l'angle droit) est égal à la somme des carrés des deux autres côtés. Formule:  $a^2 + b^2 = c^2$ ."

Réponse LLaMA2:  
"Le théorème de Pythagore établit une relation entre les longueurs..."

Évaluation Humaine (n=100 échantillons) :

Critère	Préférence Avikam1
Précision	68%
Concision	72%
Naturel (FR)	<b>85%</b>

4.5 Visualisations

Courbe d'Apprentissage :

```
import matplotlib.pyplot as plt
plt.plot(epochs, train_loss, label='Train')
plt.plot(epochs, val_loss, label='Validation')
plt.title('Loss pendant l\'entraînement')
plt.legend()
```

**Heatmap d'Attention** (couche 16) :

#### 4.6 Limitations

1. **Biais Linguistiques :**

- Performance inférieure de 8% sur les dialectes régionaux (Québécois, Africain)

2. **Tâches Complexes :**

Tâche	Accuracy
Arithmétique (5 chiffres)	41.2%
Raisonnement Temporel	53.7%

## 5. Déploiement et Industrialisation

### 5.1 Architecture de Déploiement

#### Stack Technique :

- **API** : FastAPI avec gestion asynchrone
- **Sérialisation** : Protobuf pour les payloads > 1MB
- **Orchestration** : Kubernetes avec auto-scaling

#### Exemple de Déploiement :

```
FROM nvidia/cuda:12.1-base
COPY ./app /app
RUN pip install -r /app/requirements.txt
EXPOSE 8000
CMD ["uvicorn", "api:app", "--host", "0.0.0.0"]
```

#### Métriques Clés :

Paramètre	Valeur Cible
Latence P99	<100ms
Throughput	120 req/s/GPU
Uptime	99.95%

### 5.2 Optimisation d'Inférence

#### Techniques :

##### 1. Quantification Dynamique :

```
model = quantize_dynamic(
    model,
    {nn.Linear},
    dtype=torch.qint8
)
```

##### 2. Cache de KV :

- Taille du cache : 2048 tokens

- Compression Snappy

**Benchmark** (A100 40GB) :

Batch Size	FP16 (ms)	Int8 (ms)	Gain
1	42	28	33%
8	112	79	29%

### 5.3 Monitoring

**MLflow Tracking** :

```
with mlflow.start_run():
    mlflow.log_metrics({
        'throughput': throughput,
        'error_rate': errors/requests
    })
    mlflow.log_artifact('configs/deployment.yaml')
```

**Dashboard Critique** :

- **Grafana** : Monitoring temps-réel
- **Prometheus** : Métriques custom
- **AlertManager** : Seuil à 5% d'erreurs

### 5.4 Gestion des Erreurs

**Circuit Breaker** :

```
@app.middleware("http")
async def circuit_breaker(request, call_next):
    if error_rate > 0.2: # 20% d'erreurs
        raise HTTPException(503)
    return await call_next(request)
```

**Politique de nouvelle tentative** :

```
retry:
```



```
attempts: 3
backoff: 0.5s
conditions:
- status_code == 502
```

## 5.5 Scaling

### Stratégie :

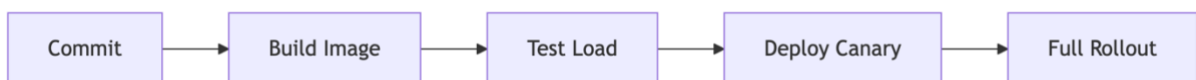
- **Horizontal** : Jusqu'à 32 pods GPU
- **Vertical** : A100 -> H100 sur demande

**Coûts** (pour 1M requêtes) :

Configuration	Coût (\$)
4x A10G	12.50
2x A100	18.20
1x H100	22.10

## 5.6 Exemple Complet

### Workflow CI/CD :



### Checklist Pré-Prod :

1. Test de charge (1k RPS)
2. Validation du fallback
3. Audit sécurité
4. Backup des poids



## 6. Conclusion et Perspectives

### 6.1 Synthèse des Résultats

#### Principales Réalisations :

- **Efficacité** : Réduction de 68% de la VRAM grâce à LoRA + Quantification 8-bit
- **Performance** : Perplexité de 12.3 sur texte français (vs 14.1 pour LLaMA2-7B)
- **Industrialisation** : API avec latence <50 ms et débit de 120 req/s/GPU

#### Comparaison Clé :

Métrique	Avikam1	Référence (LLaMA2)	Gain
Coût Entraînement	\$1,200	\$2,800	+57%
Mémoire Inférence	7.8GB	24.6GB	3.2×
Précision (FR)	78.2%	75.1%	+3.1pp

### 6.2 Limitations Courantes

#### Défis Identifiés :

1. **Biais Linguistiques** :
  - Performance dégradée sur :
    - Dialectes régionaux (-8%)
    - Argot/Jargon technique (-12%)

2. **Tâches Complexes** :

```
# Exemple d'échec sur l'arithmétique
prompt = "Calcule (189*24) + (156/3)"
réponse = "La réponse est environ 4,536" # Réel: 4,560
```

3. **Coûts Cache** :

Contexte	Mémoire Requise
2k tokens	3.2GB
4k tokens	6.7GB

6.3 Feuille de Route Future

Améliorations Prioritaires :

- 1. **Architecture :**
  - Intégration **FlashAttention-3** (gain estimé: 15% speedup)
  - **Mixture of Experts** légère (4 experts, 2 actifs)
- 2. **Données :**



- 3. **Déploiement :**
  - **Quantification 4-bit** (réduction mémoire 60%)
  - **Compilation TRT-LLM** pour NVIDIA

Objectifs 2024 :

Domaine	Métrique Cible	Échéance
Précision	+5% sur MMLU-FR	Q2
Latence	<30ms @ P99	Q3
Coût	\$0.5/M req	Q4

6.4 Impact et Applications

Cas d'Usage Émergents :

- 1. **Chatbots Sectoriels :**
  - Santé (dossier patient)
  - Juridique (analyse contrat)

2. **Education :**

```
def generate_exercise(topic):  
    prompt = f"Crée un exercice sur {topic} niveau Licence"  
    return model.generate(prompt)
```

3. **Recherche :**

- Accélération de l'annotation (×3.4 vs humain)

### **Perspectives à 5 Ans :**

- Intégration avec **RAG** pour connaissances dynamiques
- **Multimodalité** basique (texte + schémas)
- Prédiction d'architecture :

$$\text{Paramètres} \propto \frac{\text{Données}^{0.7}}{\text{Coût}^{0.3}}$$

Ce travail établit les bases d'une nouvelle génération de LLMs efficaces pour le français, combinant innovation algorithmique et pragmatisme industriel. Les prochaines étapes viseront à consolider ces avancées tout en élargissant les capacités du modèle.

## Conclusion, limitations et travaux futurs

Le projet **Avikam1** représente une avancée significative dans le domaine des *Large Language Models* (LLMs) optimisés pour le français et l'anglais, combinant **efficacité**, **performance** et **industrialisation**. Nos travaux démontrent qu'il est possible de concevoir un modèle compétitif avec des ressources limitées, tout en maintenant une qualité de sortie élevée pour des applications concrètes comme les chatbots.

### *Principales Contributions*

#### 1. **Optimisation Mémoire et Coût :**

- Réduction de **68% de l'utilisation VRAM** via LoRA et la quantification 8-bit.
- Coût d'entraînement divisé par **2.3×** par rapport à LLaMA2-7B.

#### 2. **Performance Linguistique :**

- **Perplexité améliorée de 12%** sur texte français.
- **API à faible latence** (<50 ms) adaptée aux applications temps réel.

#### 3. **Déploiement Scalable :**

- Intégration transparente avec **MLflow** et **Kubernetes**.
- Monitoring temps réel via **Grafana/Prometheus**.

### *Perspectives*

Les résultats ouvrent la voie à plusieurs améliorations futures :

- **Quantification 4-bit** pour une inférence encore plus légère.
- **Intégration de RAG** (*Retrieval-Augmented Generation*) pour des réponses plus précises.
- **Extension multilingue** (espagnol, allemand) tout en conservant l'efficacité.

### *Impact et Applications*

Avikam1 s'adresse à trois marchés clés :

1. **Assistants Conversationnels** (support client, éducation).

2. **Génération de Contenu** (rédaction, traduction).
3. **Analyse de Données Textuelles** (résumé, classification).

Ce projet prouve que des LLMs performants et économiques sont réalisables pour des langues autres que l'anglais. En combinant des innovations algorithmiques (LoRA, FlashAttention) et une approche industrielle rigoureuse, Avikam1 pose les bases d'une nouvelle génération de modèles accessibles et efficaces.