

Bien sûr, voici un plan récapitulatif que tu peux archiver pour répliquer cette architecture MoE dans un autre domaine.

Plan MoE (Phi-2 + LoRA) déjà réalisé

Préparer le socle

Répo ~/projects/soleam-photo

Venv .venv avec transformers, peft, datasets, accelerate

Script générique train\_expert\_lora.py (bf16, Trainer standard, tokenisation manuelle)

Script d'automatisation automate\_expert\_training.py

Script d'inférence moe\_router\_inference.py

Créer les datasets agents

30 fichiers agent\_\*.jsonl

Format commun : {"instruction", "response", "metadata":

{"agent", "thème", "gear", "location", "difficulty", "keywords"}}

Balises utiles pour le routage (catégorie primaire + agent précis)

Former les adaptateurs LoRA

Depuis ~/projects/soleam-photo :source .venv/bin/activate

python automate\_expert\_training.py

Résultat : models/phi2\_lora\_agent\_<nom>/final\_lora\_adapters

Configurer le routeur MoE

moe\_router\_inference.py

AGENT\_MAP et LORA\_BASE\_PATH

route\_query\_hierarchical() → catégorie (prise de vue, matériel...) puis agent (expo\_iso, portrait...)

Chargement/caching du modèle base en bf16, attachement dynamique de l'adaptateur

Fonction generate\_response(query) (prompt ### Instruction/Response)

Tester l'ensemble

Une fois 1-4 terminés :source .venv/bin/activate

python moe\_router\_inference.py

Vérifier que les logs affichent l'agent correctement routé et que la réponse exploite l'adaptateur.

Tu peux recycler ce plan pour un Smart Expert “métier X” : il suffit de reconstituer la table des agents (datasets, keywords) puis relancer les étapes 1-5.

Tu as désormais un plan reproductible pour n'importe quel autre métier :

Créer les JSONL d'agents

Par spécialité, fais 20-30 entrées instruction/response + metadata (agent, thème, difficulté, keywords...).

Pour les agents « culture/critique » de SOLÉAM, injecte les époques, artistes, maximes, etc.

Entraîner les LoRA automatiquement

Scripts : train\_expert\_lora.py (bf16, Hugging Face Trainer) et automate\_expert\_training.py.

Commande :cd ~/projects/<ton\_nouveau\_métier>

source .venv/bin/activate

python automate\_expert\_training.py

Chaque adaptateur est enregistré sous models/phi2\_lora\_agent\_<nom>/final\_lora\_adapters.

Configurer le routeur hiérarchique

moe\_router\_inference.py: mappe les mots-clés (niveau 1 catégorie, niveau 2 agent).

Ajoute une post-logique response\_clean += ... pour rappeler ton ADN de marque (ex. coupler à un agent Studio/Produit).

Tester

python moe\_router\_inference.py → vérifie que chaque question est routée vers le bon expert et qu'un adaptateur LoRA est chargé.

Ce canevas s'applique à toute nouvelle architecture « métier » : il suffit d'écrire les datasets spécialisés, de lancer l'automatisation et de définir le routage + tonalité dans l'inférence.

