

# LAD-BNet: Lag-Aware Dual-Branch Networks for Real-Time Energy Forecasting on Edge Devices

---

**Jean-Philippe LIGNIER**

*Energy Manager, Chercheur associé*

*Laboratoire de Physique et Ingénierie Mathématique pour l'Énergie, l'Environnement, et le Bâtiment*

*Université de La Réunion*

Correspondance: [jean-philippe.lignier@tangibleassets.org](mailto:jean-philippe.lignier@tangibleassets.org)

Date: 10 November 2025

---

---

## Abstract

Real-time energy forecasting on edge devices represents a major challenge for smart grid optimization and intelligent buildings. We present LAD-BNet (Lag-Aware Dual-Branch Network), an innovative neural architecture optimized for edge inference with Google Coral TPU. Our hybrid approach combines a branch dedicated to explicit exploitation of temporal lags with a Temporal Convolutional Network (TCN) featuring dilated convolutions, enabling simultaneous capture of short and long-term dependencies. Tested on real energy consumption data with 10-minute temporal resolution, LAD-BNet achieves 14.49% MAPE at 1-hour horizon with only 18ms inference time on Edge TPU, representing an 8-12 $\times$  acceleration compared to CPU. The multi-scale architecture enables predictions up to 12 hours with controlled performance degradation. Our model demonstrates a 2.39% improvement over LSTM baselines and 3.04% over pure TCN architectures, while maintaining a 180MB memory footprint suitable for embedded device constraints. These results pave the way for industrial applications in real-time energy optimization, demand management, and operational planning.

**Keywords:** energy forecasting, edge computing, neural networks, temporal convolutional networks, deep learning, Internet of Things, smart grid, edge AI

# LAD-BNet : Réseaux Neuronaux à Double Branche pour la Prévision Énergétique Temps Réel sur Dispositifs Embarqués

---

**Jean-Philippe LIGNIER**

*Energy Manager, Chercheur associé*

*Laboratoire de Physique et Ingénierie Mathématique pour l'Énergie, l'Environnement, et le Bâtiment*

*Université de La Réunion*

Correspondance: [jean-philippe.lignier@tangibleassets.org](mailto:jean-philippe.lignier@tangibleassets.org)

Date: 10 November 2025

---

---

## Résumé

La prévision énergétique en temps réel sur dispositifs embarqués représente un défi majeur pour l'optimisation des smart grids et des bâtiments intelligents. Je présente LAD-BNet (Lag-Aware Dual-Branch Network), une architecture neuronale innovante optimisée pour l'inférence edge avec Google Coral TPU. Notre approche hybride combine une branche dédiée à l'exploitation explicite des lags temporels avec un réseau convolutif temporel (TCN) à convolutions dilatées, permettant de capturer simultanément les dépendances à court et long terme. Testé sur des données réelles de consommation énergétique avec une résolution temporelle de 10 minutes, LAD-BNet atteint un MAPE de 14.49% à horizon 1h avec un temps d'inférence de seulement 18ms sur Edge TPU, soit une accélération de 8 à 12 fois par rapport au CPU. L'architecture multi-échelle permet des prévisions jusqu'à 12 heures avec une dégradation maîtrisée des performances. Notre modèle démontre une amélioration de 2.39% par rapport aux baselines LSTM et 3.04% par rapport aux architectures TCN pures, tout en maintenant une empreinte mémoire de 180MB adaptée aux contraintes des dispositifs embarqués. Ces résultats ouvrent la voie à des applications industrielles en optimisation énergétique temps réel, gestion de la demande et planification opérationnelle.

**Mots-clés:** prévision énergétique, edge computing, réseaux de neurones, TCN, deep learning, Internet des objets

# 1. Introduction

## 1.1 Contexte scientifique

La transition énergétique et le déploiement massif des énergies renouvelables nécessitent des systèmes de prévision énergétique précis et réactifs. Les smart grids modernes génèrent des volumes massifs de données temporelles haute résolution, mais leur traitement centralisé pose des problèmes de latence, de bande passante et de confidentialité. L'edge computing émerge comme paradigme prometteur, permettant le traitement local des données avec des latences minimales et une autonomie accrue (Kumar et al., 2022).

Cependant, les contraintes matérielles des dispositifs edge (mémoire limitée, puissance de calcul réduite, consommation énergétique) imposent des compromis architecturaux importants. Les modèles de deep learning traditionnels, bien que performants, sont souvent trop coûteux pour une inférence temps réel sur hardware embarqué. Les architectures optimisées pour edge doivent concilier précision prédictive, efficacité computationnelle et faible empreinte mémoire.

## 1.2 État de l'art

Les approches actuelles de prévision énergétique se divisent en plusieurs catégories. Les méthodes statistiques classiques telles que ARIMA et SARIMA offrent une interprétabilité élevée mais peinent à capturer les non-linéarités complexes des séries temporelles énergétiques (Box et al., 2015). Les réseaux de neurones récurrents, particulièrement les LSTM et GRU, dominent la littérature récente grâce à leur capacité à modéliser les dépendances temporelles longues (Hochreiter & Schmidhuber, 1997; Hewamalage et al., 2021). Cependant, leur nature séquentielle limite la parallélisation et augmente les temps d'inférence.

Les Temporal Convolutional Networks ont récemment démontré des performances compétitives avec une meilleure efficacité computationnelle (Bai et al., 2018). Les convolutions causales préservent la causalité temporelle tandis que les convolutions dilatées étendent le champ récepteur exponentiellement (Bai et al., 2018). Néanmoins, les TCN purs négligent parfois l'exploitation explicite des lags temporels, pourtant fondamentaux en prévision énergétique où les autocorrelations sont fortes.

Les architectures Transformer adaptées aux séries temporelles ont également émergé récemment, avec des modèles spécialisés comme Informer (Zhou et al., 2021) et Autoformer (Wu et al., 2021) démontrant des performances compétitives. Les approches d'expansion par bases neurales comme N-BEATS (Oreshkin et al., 2019) offrent une alternative interprétable aux méthodes boîte noire. L'accélération hardware via TPU et NPU a transformé le paysage de l'edge AI. Google Coral offre des performances d'inférence remarquables pour les modèles TensorFlow Lite quantifiés, mais impose des contraintes architecturales strictes (Google Coral Team, 2025). La conception de modèles exploitant pleinement ces accélérateurs tout en maintenant la précision prédictive reste un défi ouvert.

### **1.3 Problématique et contributions**

Cette recherche répond à la question suivante : comment concevoir une architecture neuronale pour la prévision énergétique temps réel qui soit simultanément précise, rapide et compatible avec les contraintes des dispositifs edge ? Nous identifions trois limitations majeures des approches existantes : l'exploitation insuffisante des lags dans les TCN, l'absence d'architectures hybrides optimisées pour edge TPU, et le manque de solutions déployables en production industrielle.

Mes contributions principales sont les suivantes :

- Proposition de LAD-BNet, une architecture dual-branch innovante combinant exploitation explicite des lags et convolutions temporelles dilatées
- Démonstration d'une amélioration significative des performances prédictives sur des données réelles de consommation énergétique
- Validation de l'efficacité computationnelle sur Google Coral TPU avec des temps d'inférence inférieurs à 20ms
- Fourniture d'une solution complète production-ready incluant monitoring, alertes et dashboard temps réel

### **1.4 Organisation de l'article**

Cet article s'organise comme suit. La section 2 détaille l'architecture LAD-BNet et les choix méthodologiques. La section 3 présente les données expérimentales et le protocole d'évaluation. La section 4 expose les résultats quantitatifs et qualitatifs. La section 5 discute les implications, limitations et perspectives. La section 6 conclut sur les contributions et applications industrielles.

## 2. Matériel et Méthodes

### 2.1 Architecture LAD-BNet

#### 2.1.1 Philosophie de conception

LAD-BNet repose sur l'hypothèse que les séries temporelles énergétiques contiennent à la fois des dépendances explicites (lags) et des patterns implicites (motifs récurrents). Une architecture dual-branch permet de spécialiser chaque branche dans l'extraction d'un type de features, puis de fusionner ces représentations complémentaires.

#### 2.1.2 Architecture globale

Le modèle accepte en entrée une séquence de 144 pas temporels (24 heures à résolution 10 minutes) avec 27 features par pas. L'architecture se décompose en trois modules principaux :

Branch 1 (Lag Branch) : Cette branche exploite les 24 derniers pas temporels (4 heures) via un réseau dense profond. Les lags temporels ( $kW\_lag\_6$ ,  $kW\_lag\_12$ ,  $kW\_lag\_24$ ,  $kW\_lag\_72$ ,  $kW\_lag\_144$ ) présentent des corrélations fortes avec la consommation future. L'architecture séquentielle comprend un aplatissement de la séquence vers un vecteur de dimension 648, suivi de deux couches denses de dimensions 256 et 128 avec normalisation batch et dropout 0.1. Cette configuration permet d'encoder efficacement la mémoire explicite du système.

Branch 2 (TCN Branch) : Cette branche traite la séquence complète via des convolutions temporelles. Deux convolutions causales de 64 filtres avec noyaux de taille 3 capturent les patterns locaux. Une convolution dilatée de 128 filtres avec facteur de dilatation 2 étend le champ récepteur à l'ensemble de la séquence. Le dual pooling (average et max) extrait les features les plus saillantes. Cette architecture permet de capturer les dépendances temporelles à différentes échelles sans augmenter exponentiellement les paramètres.

Module de fusion : Les représentations des deux branches sont concaténées puis projetées via des couches denses de dimensions 256 et 128 avec dropout. La couche de sortie génère 72 prédictions correspondant à un horizon de 12 heures.

#### 2.1.3 Justification des choix architecturaux

Le dual-branch design est motivé par l'analyse de corrélation des features. Les lags présentent un coefficient de détermination  $R^2$  supérieur à 0.85 avec la consommation future, justifiant une branche dédiée. Les convolutions dilatées offrent un champ récepteur exponentiel : avec un noyau de taille 3 et dilatation 2, le champ récepteur atteint 7 pas temporels en une seule couche (Bai et al., 2018).

Le dropout uniforme à 0.1 est un compromis entre régularisation et capacité d'apprentissage. Des expériences préliminaires ont montré qu'un dropout supérieur à 0.2 dégrade les performances. La normalisation batch stabilise l'entraînement et accélère la convergence. Le dual pooling capture simultanément la tendance moyenne et les extrema, information cruciale pour la prévision de pics énergétiques.

## 2.2 Engineering des features

### 2.2.1 Features temporelles et contextuelles

La construction du vecteur de features repose sur l'extraction d'information multi-échelle. Les features temporelles cycliques (hour\_sin/cos, month\_sin/cos, dayofweek\_sin/cos) encodent la périodicité naturelle via transformation trigonométrique, évitant la discontinuité artificielle ( $23h \rightarrow 0h$ ) (Lim & Zohren, 2021). Les features contextuelles binaires (weekend, is\_holiday, is\_business\_hours, is\_night, is\_morning\_peak, is\_evening\_peak) capturent les régimes opérationnels distincts du bâtiment.

### 2.2.2 Features de lags et statistiques roulantes

Les lags temporels à horizons multiples (1h, 2h, 4h, 12h, 24h) fournissent une mémoire explicite du système. Les statistiques roulantes (moyenne, écart-type, maximum, minimum) sur fenêtres glissantes (1h, 2h, 4h) caractérisent la volatilité et les tendances locales. Ces features agrégées réduisent la sensibilité au bruit de mesure.

### 2.2.3 Features météorologiques et interactions

Les features météorologiques (température de bulle sec DBT, humidité relative RH) influencent fortement la consommation via la charge de climatisation. La feature d'interaction température-humidité ( $DBT \times RH / 100$ ) capture les effets non-linéaires sur le confort thermique et donc sur la demande énergétique.

## 2.3 Préparation des données

### 2.3.1 Normalisation

Toutes les features sont normalisées via MinMaxScaler dans l'intervalle [0,1]. Cette normalisation garantit la stabilité numérique et accélère la convergence du gradient descent. Les paramètres du scaler (minimum et maximum par feature) sont sauvegardés pour la dénormalisation post-inférence.

### 2.3.2 Fenêtrage temporel

Les données sont transformées en séquences via une fenêtre glissante. Chaque échantillon d'entraînement comprend 144 pas temporels d'entrée (24h) et 72 pas de sortie (12h). Un stride de 1 génère des séquences hautement corrélées, ce qui favorise la généralisation mais augmente le dataset apparent. En production, la fenêtre glisse avec l'arrivée de nouvelles mesures toutes les 10 minutes.

## 2.4 Entraînement et optimisation

### 2.4.1 Configuration d'entraînement

L'entraînement utilise l'optimiseur Adam avec learning rate 0.0005, choisi après grid search sur [0.0001, 0.001]. La fonction de perte est l'erreur quadratique moyenne (MSE), adaptée à la régression temporelle. La batch size de 16 équilibre stabilité du gradient et vitesse d'entraînement. L'entraînement s'effectue sur 400 époques avec early stopping basé sur la validation loss (patience 50 époques).

## 2.4.2 Régularisation

Outre le dropout, j'applique la normalisation batch après chaque couche dense et convulsive. Cette double régularisation prévient l'overfitting tout en accélérant la convergence. Les poids sont initialisés via glorot\_uniform, optimale pour les activations ReLU.

## 2.4.3 Split train/validation/test

Les données sont divisées chronologiquement : 70% entraînement, 15% validation, 15% test. Le split chronologique préserve la distribution temporelle et évite le data leakage. La validation guide la sélection des hyperparamètres, tandis que le test évalue la généralisation finale.

## 2.5 Déploiement sur Edge TPU

### 2.5.1 Conversion TensorFlow Lite

Le modèle Keras entraîné est converti en TensorFlow Lite avec quantification full integer int8 (Jacob et al., 2018). Cette quantification réduit la taille du modèle de 75% et accélère l'inférence sur TPU (Han et al., 2015). Le processus de quantification utilise un dataset représentatif pour calibrer les échelles de quantification, minimisant la perte de précision.

### 2.5.2 Compilation Edge TPU

Le modèle TFLite est compilé pour Edge TPU via le compilateur edgetpu\_compiler. Toutes les opérations sont mappées sur le TPU sauf quelques opérations post-processing mineures qui s'exécutent sur CPU. Le taux de compilation Edge TPU atteint 95%, garantissant une accélération maximale.

### 2.5.3 Infrastructure de déploiement

Le système de production s'exécute sur Raspberry Pi 4B (4GB RAM) équipé de Google Coral USB Accelerator. Le runtime PyCoral gère l'interface avec le TPU. Un service systemd garantit le redémarrage automatique en cas d'erreur. Le dashboard Streamlit offre une visualisation temps réel accessible via navigateur web.

## 2.6 Métriques d'évaluation

### 2.6.1 Précision prédictive

La métrique principale est le Mean Absolute Percentage Error (MAPE), défini comme :

$$\text{MAPE} = (100/n) \times \sum |y_{\text{réel}} - y_{\text{prédict}}| / |y_{\text{réel}}|$$

Le MAPE offre une interprétation intuitive en pourcentage et est robuste aux changements d'échelle. Je calcule également le coefficient de détermination  $R^2$  pour quantifier la variance expliquée.

### 2.6.2 Performance d'inférence

Je mesure le temps d'inférence de bout en bout, incluant le prétraitement des features, l'inférence TPU et la dénormalisation. Les percentiles P50, P95 et P99 caractérisent la distribution des latences. Le throughput en prédictions par seconde quantifie la capacité de traitement.

### **2.6.3 Consommation de ressources**

Je monitore l'utilisation mémoire (RAM), le taux d'occupation CPU, et la température du processeur. Ces métriques valident la viabilité du déploiement continu sur dispositif embarqué.

## **2.7 Protocole expérimental**

### **2.7.1 Baseline et comparaisons**

Je compare LAD-BNet à trois baselines : un LSTM bidirectionnel (2 couches de 128 unités), un TCN pur (architecture V6.0 sans branch lags), et une approche statistique SARIMA. Les baselines sont entraînées sur les mêmes données avec optimisation des hyperparamètres.

### **2.7.2 Validation multi-horizon**

Les performances sont évaluées à 5 horizons de prévision : 1h, 2h, 4h, 8h et 12h. Cette évaluation multi-horizon révèle la dégradation des performances avec l'horizon, caractéristique fondamentale de tout modèle prédictif.

### **2.7.3 Reproductibilité**

Tous les seeds aléatoires (numpy, tensorflow) sont fixés pour garantir la reproductibilité. Le code source, les configurations et les données (anonymisées) sont disponibles pour réPLICATION indépendante.

### 3. Résultats

#### 3.1 Performances prédictives

##### 3.1.1 MAPE multi-horizon

Le tableau 1 présente les performances de LAD-BNet à différents horizons de prévision.

Tableau 1 : Performances du modèle LAD-BNet V7.1

Horizon	Pas temporels	MAPE (%)	R <sup>2</sup>	Objectif
1h	6	14.49	0.869	< 15%
2h	12	14.82	0.861	< 16%
4h	24	15.31	0.848	< 17%
8h	48	15.68	0.836	< 18%
12h	72	16.31	0.826	< 19%

LAD-BNet atteint un MAPE de 14.49% à horizon 1h, surpassant l'objectif de 15%. La dégradation progressive avec l'horizon ( $\Delta +1.82$  points de 1h à 12h) démontre la robustesse du modèle. Le coefficient R<sup>2</sup> reste supérieur à 0.82 même à 12h, indiquant une forte capacité explicative.

##### 3.1.2 Comparaison avec baselines

Le tableau 2 compare LAD-BNet aux approches concurrentes.

Tableau 2 : Comparaison avec baselines à horizon 1h

Modèle	MAPE (%)	R <sup>2</sup>	Amélioration vs Baseline
SARIMA	22.15	0.621	-
LSTM Bidirectionnel	16.88	0.805	Baseline
TCN V6.0 (pur)	17.92	0.782	-6.16%
LAD-BNet V7.1	14.49	0.869	+14.16%

LAD-BNet améliore le MAPE de 2.39 points par rapport au LSTM baseline (Salinas et al., 2020) et de 3.43 points par rapport au TCN pur. L'architecture dual-branch démontre une supériorité claire sur les approches mono-branche.

##### 3.1.3 Analyse temporelle des erreurs

La figure 1 (non reproduite) illustre l'évolution du MAPE sur une semaine test. Les erreurs augmentent légèrement durant les weekends (MAPE ~15.8%) par rapport aux jours ouvrés (MAPE ~14.1%), reflétant des patterns de consommation plus variables. Les heures de

transition (7h-9h, 17h-20h) présentent des erreurs légèrement supérieures dues aux changements rapides de charge.

## 3.2 Performances d'inférence

### 3.2.1 Temps d'inférence

Le tableau 3 détaille les performances d'inférence sur différents dispositifs.

**Tableau 3 : Temps d'inférence et accélération hardware**

Dispositif	Temps moyen	P50	P95	P99	Accélération
<b>CPU Raspberry Pi 4B</b>	150 ms	145 ms	180 ms	210 ms	1x
<b>Edge TPU (standard)</b>	18 ms	17 ms	23 ms	28 ms	8.3x
<b>Edge TPU (max)</b>	12 ms	11 ms	16 ms	20 ms	12.5x

L'accélération Edge TPU varie de 8.3x (mode standard) à 12.5x (mode max). La latence P99 de 28ms en mode standard garantit la contrainte temps réel de 100ms. Le mode max réduit encore la latence mais augmente la température du TPU de 10-15°C.

### 3.2.2 Throughput

Le système atteint un throughput de 55 prédictions par seconde (55 Hz) en mode continu. Cela correspond à 198,000 prédictions par heure, largement supérieur au besoin opérationnel d'une prédiction toutes les 10 minutes. Ce surdimensionnement permet le traitement de multiples flux en parallèle ou l'exécution d'ensemble de modèles.

### 3.2.3 Consommation de ressources

Le tableau 4 quantifie l'empreinte système.

**Tableau 4 : Utilisation des ressources système**

Métrique	Idle	Inférence	Maximum observé
<b>RAM</b>	150 MB	180 MB	195 MB
<b>CPU</b>	< 5%	15-25%	32%
<b>Température CPU</b>	45°C	52°C	58°C
<b>Température TPU</b>	38°C	48°C	55°C

L'empreinte mémoire reste inférieure à 200MB, compatible avec les 4GB du Raspberry Pi 4B. Le CPU demeure largement disponible pour les tâches annexes (monitoring, dashboard). Les températures restent dans les limites sécuritaires sans nécessiter de refroidissement actif.

### 3.3 Scalabilité et multi-échelle

#### 3.3.1 Scalabilité horizontale

Des tests sur multiples dispositifs Raspberry Pi démontrent une scalabilité linéaire. Un cluster de 4 Pi équipés de Coral atteint 220 prédictions/seconde avec une latence moyenne maintenue à 18ms. Cette scalabilité horizontale permet le déploiement dans des environnements multi-bâtiments.

#### 3.3.2 Adaptabilité multi-échelle

Le modèle s'adapte à différentes résolutions temporelles via ré entraînement. Des expériences préliminaires à résolution 5 minutes (au lieu de 10) maintiennent un MAPE de 15.2% à horizon 1h. À résolution 30 minutes, le MAPE descend à 13.1%, confirmant que la précision augmente avec l'agrégation temporelle.

#### 3.3.3 Transfert learning

Des tests de transfert learning montrent que le modèle pré-entraîné sur un bâtiment converge 3x plus rapidement sur un nouveau bâtiment (50 époques vs 150). Le fine-tuning de la couche de sortie uniquement atteint 92% des performances du modèle complet ré entraîné, avec seulement 5% du coût computationnel.

### 3.4 Analyse d'ablation

#### 3.4.1 Contribution des branches

Le tableau 5 quantifie l'apport de chaque branche.

Tableau 5 : Étude d'ablation des composantes

Configuration	MAPE 1h (%)	Δ vs Complet	Temps (ms)
LAD-BNet complet	14.49	0	18
Lag Branch seule	16.22	+1.73	12
TCN Branch seule	17.92	+3.43	15
Sans conv. dilatées	15.87	+1.38	16
Sans dual pooling	15.12	+0.63	17

L'ablation révèle que les deux branches apportent des contributions complémentaires. La TCN Branch seule retrouve les performances du TCN V6.0, confirmant que l'amélioration provient de la fusion avec la Lag Branch. Les convolutions dilatées et le dual pooling contribuent marginalement mais significativement.

#### 3.4.2 Impact des features

Une analyse SHAP révèle que les lags (kW\_lag\_6, kW\_lag\_12) contribuent à 35% de la prédiction, les features temporelles cycliques à 25%, les statistiques roulantes à 20%, les features contextuelles à 15% et les features météorologiques à 5%. Cette distribution justifie l'architecture dual-branch priorisant les lags.

## 3.5 Robustesse et fiabilité

### 3.5.1 Robustesse aux données manquantes

Des simulations avec 5%, 10% et 20% de données manquantes (interpolation linéaire) augmentent le MAPE de respectivement 0.3%, 0.8% et 2.1%. Le modèle dégrade gracieusement jusqu'à 10% de données manquantes, au-delà la dégradation s'accélère.

### 3.5.2 Stabilité long-terme

Un déploiement continu sur 30 jours montre une dérive du MAPE de +0.5% par rapport aux performances initiales, attribuable à l'évolution des patterns de consommation. Un ré entraînement hebdomadaire maintient les performances à  $\pm 0.2\%$  du niveau initial.

### 3.5.3 Fiabilité système

Sur 10,000 prédictions en conditions réelles, le taux de succès atteint 99.92% avec 8 échecs attribuables à des déconnexions USB du Coral. Le système systemd redémarre automatiquement le service en moins de 5 secondes après un crash, garantissant une disponibilité de 99.98%.

## 4. Discussion

### 4.1 Interprétation des résultats

#### 4.1.1 Supériorité de l'architecture dual-branch

Les résultats démontrent la supériorité de LAD-BNet sur les architectures mono-branche. Cette performance s'explique par l'exploitation de deux types de dépendances temporelles. La Lag Branch capture les autocorrélations fortes à court terme, phénomène dominant dans les séries énergétiques où la consommation actuelle prédit fortement la consommation future immédiate (Qin et al., 2017). La TCN Branch identifie les patterns récurrents (cycles quotidiens, hebdomadaires) via les convolutions qui apprennent des filtres temporels adaptatifs.

La fusion tardive des représentations permet une spécialisation optimale de chaque branche. Des expériences avec fusion précoce (concaténation avant projection) dégradent les performances de 1.2 points de MAPE, confirmant l'importance de l'extraction de features hiérarchiques.

#### 4.1.2 Efficacité des convolutions dilatées

Les convolutions dilatées étendent le champ récepteur sans augmenter le nombre de paramètres. Avec un noyau de taille 3 et dilatation 2, le champ récepteur effectif couvre 7 pas temporels en une seule couche. Pour capturer les 144 pas d'entrée, seulement 3-4 couches dilatées suffisent, contre 8-10 couches convolutives standard. Cette efficacité architecturale réduit la profondeur du réseau et accélère l'inférence.

La dilatation préserve également la résolution temporelle en évitant le pooling agressif. Cette propriété est cruciale pour la prévision énergétique où les variations haute fréquence (pics de consommation) contiennent une information prédictive importante.

#### 4.1.3 Impact de la quantification

La quantification int8 réduit la taille du modèle de 12MB (float32) à 3MB (int8) sans dégradation significative des performances. Le MAPE augmente de seulement 0.15% après quantification, une perte négligeable au regard du gain en vitesse d'inférence (8-12x). Cette robustesse s'explique par la normalisation préalable des features dans [0,1] et la nature bornée des prédictions.

## 4.2 Positionnement par rapport aux travaux antérieurs

### 4.2.1 Comparaison avec la littérature

Le MAPE de 14.49% à 1h se compare favorablement à la littérature récente sur la prévision énergétique. Zhang et al. (2020) rapportent un MAPE de 16.2% avec LSTM sur données similaires. Li et al. (2021) atteignent 15.8% avec Transformer, mais avec un temps d'inférence de 200ms sur GPU. Notre approche offre un meilleur compromis précision-latence adapté au edge computing.

Les travaux spécifiques sur edge AI pour l'énergie sont rares. Kumar et al. (2022) déplient un LSTM quantifié sur Raspberry Pi avec un MAPE de 17.5% et une latence de 80ms sans TPU. L'accélération via Coral représente une avancée significative pour les applications temps réel.

## **4.2.2 Originalité de l'approche**

La contribution principale réside dans l'architecture dual-branch optimisée pour Edge TPU. À ma connaissance, aucun travail antérieur ne combine explicitement lags et TCN dans une architecture unifiée pour l'edge computing énergétique. La conception guidée par les contraintes hardware (quantification int8, opérations compatibles TPU) différencie notre approche des modèles académiques non déployables.

## **4.3 Applications industrielles**

### **4.3.1 Optimisation énergétique en temps réel**

LAD-BNet permet l'implémentation de stratégies de demand response avec latence minimale, contribuant à l'optimisation des smart grids (Zhang et al., 2020; Kumar et al., 2022) et des bâtiments intelligents (Li et al., 2021). La prévision des pics de consommation 1-2h à l'avance autorise le délestage préventif ou l'activation de stockage énergétique. Un déploiement pilote dans un bâtiment tertiaire a réduit la consommation durant les heures de pointe de 8.3% via ajustement automatique de la climatisation basé sur les prévisions.

### **4.3.2 Gestion intelligente des microgrids**

Dans les microgrids avec production renouvelable (Zhang et al., 2020; Kumar et al., 2022), la prévision de consommation complète la prévision de production pour optimiser l'équilibrage offre-demande.. LAD-BNet déployé sur des contrôleurs edge locaux évite la centralisation et réduit les délais de communication. Une installation expérimentale sur microgrid insulaire a amélioré l'efficacité d'utilisation du stockage de 12% grâce aux prédictions locales.

### **4.3.3 Maintenance prédictive**

Les anomalies dans les prédictions signalent des dysfonctionnements potentiels. Une consommation réelle systématiquement supérieure aux prédictions peut indiquer une dégradation d'isolation ou une défaillance de système HVAC. Le monitoring continu facilite la détection précoce de dérives avant panne majeure.

### **4.3.4 Participation aux marchés énergétiques**

Les agrégateurs d'effacement utilisent les prévisions de consommation pour optimiser leurs offres sur les marchés day-ahead et intraday. LAD-BNet fournit les prévisions multi-horizons (1h-12h) nécessaires à ces applications. La latence sub-30ms permet le recalcul fréquent des prévisions pour réactivité maximale aux conditions changeantes.

### **4.3.5 Intégration dans les systèmes de gestion de bâtiment (BMS)**

Le déploiement sur Raspberry Pi facilite l'intégration dans l'infrastructure existante des bâtiments intelligents. LAD-BNet communique via API REST avec les BMS pour fournir des prédictions consommées par les contrôleurs HVAC, d'éclairage et de sécurité. Un prototype d'intégration avec BACnet démontre l'interopérabilité avec les standards industriels.

## **4.4 Limites de l'étude**

### **4.4.1 Généralisation aux différents types de bâtiments**

Les expériences se concentrent sur des données d'un bâtiment tertiaire universitaire. La généralisation à d'autres typologies (résidentiel, industriel, commercial) nécessite validation.

Les patterns de consommation varient significativement entre secteurs, et un réentraînement adapté peut être nécessaire. Nos tests de transfert learning suggèrent néanmoins que l'architecture se transpose avec réentraînement limité.

#### 4.4.2 Dépendance aux prévisions météorologiques

Le modèle utilise des mesures météorologiques passées mais pas de prévisions météo futures. L'intégration de prévisions météo améliorerait probablement la précision à horizons longs ( $> 4\text{h}$ ). Cette limitation provient de l'indisponibilité locale des prévisions météo, mais une intégration via API est envisageable.

#### 4.4.3 Horizon de prévision limité

LAD-BNet prédit jusqu'à 12h, insuffisant pour certaines applications (planification day-ahead des marchés énergétiques). L'extension à 24-48h dégrade significativement les performances (MAPE  $> 20\%$ ). Une approche multi-modèle avec réajustement fréquent pourrait pallier cette limitation.

#### 4.4.4 Coût hardware

Le Coral USB Accelerator (70-80€) et le Raspberry Pi 4B 8GB (80-90€) représentent un investissement initial d'environ 150-170€ par nœud. Pour des déploiements massifs, ce coût peut être prohibitif. Une analyse coût-bénéfice intégrant les économies énergétiques est nécessaire pour justifier l'investissement.

### 4.5 Perspectives de recherche

Les développements futurs de LAD-BNet s'articulent autour de plusieurs axes :

#### 4.5.1 Architecture attention-based

L'intégration de mécanismes d'attention (self-attention, cross-attention) pourrait améliorer la modélisation des dépendances long-terme (Vaswani et al., 2017; Li et al., 2021). Les Transformers spécialisées pour séries temporelles (Wen et al., 2022; Zhou et al., 2021) démontrent des performances remarquables mais leur coût computationnel reste élevé. Une architecture hybride TCN-Attention optimisée pour edge constitue une piste prometteuse.

#### 4.5.2 Apprentissage multi-tâches

L'extension à l'apprentissage multi-tâches (prévision conjointe de consommation électrique, thermique, et production renouvelable) pourrait exploiter les corrélations entre tâches. Un encodeur partagé avec têtes de prédiction spécialisées réduirait les paramètres totaux tout en améliorant potentiellement chaque tâche.

#### 4.5.3 Quantification adaptative

La quantification actuelle (int8 uniforme) est sous-optimale. Certaines couches critiques bénéficieraient d'une précision supérieure (int16) tandis que d'autres tolèrent int4. Une quantification mixte guidée par l'analyse de sensibilité optimiserait le compromis précision-vitesse.

#### 4.5.4 Apprentissage fédéré

Pour des déploiements multi-sites, l'apprentissage fédéré permettrait l'amélioration continue du modèle sans centralisation des données. Chaque nœud edge entraînerait localement puis

partagerait les gradients agrégés. Cette approche préserve la confidentialité tout en capitalisant sur la diversité des sites.

#### **4.5.5 Prévision probabiliste**

Les prédictions actuelles sont ponctuelles (point forecasts). L'extension à des prévisions probabilistes (intervalles de confiance, quantiles) fournirait une information d'incertitude cruciale pour la prise de décision en conditions risquées (Gasthaus et al., 2019; Salinas et al., 2020). Des architectures bayésiennes ou des ensembles de modèles sont envisageables.

### **4.6 Implications pour l'edge AI**

#### **4.6.1 Design guidé par les contraintes hardware**

La méthodologie démontre l'importance du co-design modèle-hardware. La conception de LAD-BNet intègre dès l'origine les contraintes Edge TPU (opérations supportées, quantification). Cette approche contraste avec l'entraînement cloud puis compression post-hoc, souvent sous-optimale.

#### **4.6.2 Compromis précision-efficacité**

Les résultats quantifient le compromis précision-latence-mémoire fondamental en edge AI. LAD-BNet sacrifie 0.5-1% de précision (vs modèles cloud) pour gagner 10x en vitesse et réduire la mémoire de 75%. Ce compromis est acceptable pour la majorité des applications industrielles temps réel.

#### **4.6.3 Viabilité du edge computing pour l'énergie**

L'étude valide la viabilité technique et économique du edge computing pour la prévision énergétique. Les performances démontrées (MAPE < 15%, latence < 20ms, coût < 200€) ouvrent la voie à un déploiement massif dans les smart buildings et smart grids distribués.

## 5. Conclusion

### 5.1 Synthèse des contributions

Cette recherche présente LAD-BNet, une architecture neuronale innovante pour la prévision énergétique temps réel sur dispositifs edge. Mes contributions principales sont les suivantes.

Premièrement, je propose une architecture dual-branch combinant exploitation explicite des lags temporels et convolutions temporelles dilatées. Cette conception hybride capture simultanément les dépendances à court et long terme, surpassant les approches mono-branche de 2.39 à 3.43 points de MAPE.

Deuxièmement, je démontre l'efficacité d'un déploiement optimisé sur Google Coral Edge TPU. L'accélération hardware de 8 à 12 fois réduit la latence à 18ms tout en maintenant une empreinte mémoire de 180MB, compatible avec les dispositifs embarqués.

Troisièmement, je valide les performances sur des données réelles de consommation énergétique. Le MAPE de 14.49% à horizon 1h et la dégradation maîtrisée jusqu'à 12h (16.31%) positionnent LAD-BNet parmi les approches state-of-the-art pour l'edge computing énergétique.

Quatrièmement, je fournis une solution complète production-ready incluant pipeline de données, monitoring continu, système d'alertes et dashboard de visualisation. Cette approche système facilite l'adoption industrielle immédiate.

### 5.2 Retombées pratiques

Les applications industrielles de LAD-BNet sont multiples et immédiates. L'optimisation énergétique en temps réel permet des économies de 5-10% sur les coûts d'énergie via demand response et effacement. La gestion intelligente des microgrids améliore l'intégration des renouvelables et réduit le besoin de stockage coûteux. La maintenance prédictive basée sur la détection d'anomalies prévient les pannes majeures. La participation optimisée aux marchés énergétiques augmente les revenus des agrégateurs.

Au-delà de l'énergie, LAD-BNet démontre la viabilité du edge AI pour les séries temporelles. L'architecture se transpose potentiellement à d'autres domaines tels que la prévision de trafic réseau, de flux de production industrielle, ou de demande de services cloud.

### 5.3 Retombées théoriques

Sur le plan théorique, cette recherche contribue à la compréhension des architectures neurales hybrides pour séries temporelles. La démonstration que l'exploitation explicite de features (lags) combinée à l'apprentissage implicite (convolutions) surpassé les approches purement end-to-end remet en question le paradigme dominant du deep learning "black-box".

L'étude quantitative du compromis précision-efficacité en edge computing enrichit la littérature sur le co-design modèle-hardware. Les résultats suggèrent que des architectures spécifiquement conçues pour l'edge peuvent égaler voire surpasser les modèles cloud sur des métriques pondérées (précision × vitesse × efficacité énergétique).

### 5.4 Perspectives d'évolution

Les développements futurs de LAD-BNet s'articulent autour de trois axes. L'amélioration de l'architecture via intégration de mécanismes d'attention et apprentissage multi-tâches promet des gains de précision additionnels. L'extension à des prévisions probabilistes et

l'apprentissage fédéré répondront aux besoins d'applications critiques. Le déploiement sur hardware de nouvelle génération (Coral Max, Intel Movidius Myriad X, Nvidia Jetson Nano) explorera les limites de performance de l'edge AI.

Au-delà des améliorations techniques, la validation sur datasets publics (UCI, IHEPC) et la comparaison systématique avec méthodes state-of-the-art renforceront la robustesse scientifique. Un déploiement à large échelle dans un contexte smart city quantifiera l'impact économique et environnemental à l'échelle urbaine.

## 5.5 Conclusion finale

LAD-BNet démontre que la prévision énergétique précise et temps réel sur dispositifs edge est techniquement viable et économiquement attractive. L'architecture dual-branch, l'optimisation Edge TPU et la solution système complète convergent pour offrir une solution industrialisable immédiatement. Les performances quantitatives (MAPE 14.49%, latence 18ms, empreinte 180MB) positionnent LAD-BNet comme référence pour l'edge AI énergétique.

L'impact potentiel est considérable. Le déploiement massif de LAD-BNet dans les smart buildings pourrait réduire la consommation énergétique globale de 5-10%, contribuant significativement aux objectifs de neutralité carbone. La décentralisation du traitement vers l'edge améliore la résilience des systèmes énergétiques face aux cyberattaques et pannes centralisées.

En conclusion, LAD-BNet représente une avancée significative dans le domaine émergent de l'edge AI pour l'énergie. L'approche proposée, combinant rigueur scientifique et pragmatisme industriel, trace la voie vers des systèmes énergétiques intelligents, autonomes et efficaces. Les développements futurs promettent d'amplifier cet impact, rapprochant la vision de smart grids véritablement distribués et réactifs.

## 6. Remerciements

Je remercie le Laboratoire de Physique et Ingénierie Mathématique pour l'Énergie, l'Environnement et le Bâtiment (PIMENT) de l'Université de La Réunion pour le support institutionnel qui m'a été offert en tant que chercheur associé et pour la mise à disposition des données de consommation électrique et métrologiques qui ont permis la conception de ce modèle. J'exprime ma gratitude particulière à Mathieu David, directeur du laboratoire PIMENT, pour sa bienveillance et ses observations toujours pertinentes.

Cette recherche a bénéficié des discussions enrichissantes avec les membres du laboratoire et d'échanges constructifs avec des industriels du secteur énergétique réunionnais. Mes remerciements s'adressent également à Khalil Drira, Directeur de Recherche au CNRS (équipe SARA – Services et Architectures pour les Réseaux Avancés, LAAS-CNRS, Université de Toulouse), pour ses conseils toujours judicieux lors de nos échanges.

L'intégralité de l'infrastructure matérielle utilisée dans cette recherche a été autofinancée. Le système de développement (station de travail équipée d'un processeur Intel Core i5-13400F, 64 GB RAM, GPU NVIDIA RTX 4070) a été acquis d'occasion, tandis que le système de production edge (Raspberry Pi 4B 8GB avec Google Coral USB Accelerator) a été acquis neuf. Cette approche démontre la faisabilité de recherches en edge AI avec des ressources limitées et un impact environnemental maîtrisé.

Je remercie la communauté open-source, en particulier les contributeurs de Python, PyTorch, TensorFlow, scikit-learn et PyCoral, dont les outils ont été essentiels à ce travail.

Enfin, je remercie mon épouse Valérie pour son soutien indéfectible et son engagement partagé à développer des solutions performantes à faible coût et à faible impact environnemental.

## 7. Références bibliographiques

- [1] Bai, S., Kolter, J. Z., & Koltun, V. (2018). An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. arXiv preprint arXiv:1803.01271.
- [2] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. Neural computation, 9(8), 1735-1780.
- [3] Zhang, Y., Wang, J., & Chen, X. (2020). Deep learning for energy consumption prediction in smart buildings. IEEE Transactions on Industrial Informatics, 16(6), 4145-4155.
- [4] Li, W., Zhao, H., & Zhang, L. (2021). Transformer-based energy forecasting with multi-scale attention. Applied Energy, 295, 117045.
- [5] Kumar, S., Singh, R., & Patel, M. (2022). Edge computing for real-time energy management: A Raspberry Pi implementation. Journal of Ambient Intelligence and Smart Environments, 14(2), 145-162.
- [6] Box, G. E., Jenkins, G. M., Reinsel, G. C., & Ljung, G. M. (2015). Time series analysis: forecasting and control (5th ed.). John Wiley & Sons.
- [7] Hewamalage, H., Bergmeir, C., & Bandara, K. (2021). Recurrent neural networks for time series forecasting: Current status and future directions. International Journal of Forecasting, 37(1), 388-427.
- [8] Oreshkin, B. N., Carpov, D., Chapados, N., & Bengio, Y. (2019). N-BEATS: Neural basis expansion analysis for interpretable time series forecasting. arXiv preprint arXiv:1905.10437.
- [9] Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). Attention is all you need. Advances in neural information processing systems, 30.
- [10] Jacob, B., Kligys, S., Chen, B., et al. (2018). Quantization and training of neural networks for efficient integer-arithmetic-only inference. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 2704-2713).
- [11] Han, S., Mao, H., & Dally, W. J. (2015). Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. arXiv preprint arXiv:1510.00149.
- [12] Ahmad, T., Chen, H., Guo, Y., & Wang, J. (2018). A comprehensive overview on the data driven and large scale based approaches for forecasting of building energy demand: A review. Energy and Buildings, 165, 301-320.
- [13] Raza, M. Q., & Khosravi, A. (2015). A review on artificial intelligence based load demand forecasting techniques for smart grid and buildings. Renewable and Sustainable Energy Reviews, 50, 1352-1372.
- [14] Mocanu, E., Nguyen, P. H., Gibescu, M., & Kling, W. L. (2016). Deep learning for estimating building energy consumption. Sustainable Energy, Grids and Networks, 6, 91-99.
- [15] Marino, D. L., Amarasinghe, K., & Manic, M. (2016). Building energy load forecasting using deep neural networks. In IECON 2016-42nd Annual Conference of the IEEE Industrial Electronics Society (pp. 7046-7051).

- [16] Kong, W., Dong, Z. Y., Jia, Y., et al. (2019). Short-term residential load forecasting based on LSTM recurrent neural network. *IEEE Transactions on Smart Grid*, 10(1), 841-851.
- [17] Bouktif, S., Fiaz, A., Ouni, A., & Serhani, M. A. (2018). Optimal deep learning LSTM model for electric load forecasting using feature selection and genetic algorithm: Comparison with machine learning approaches. *Energies*, 11(7), 1636.
- [18] Wang, H., Lei, Z., Zhang, X., Zhou, B., & Peng, J. (2019). A review of deep learning for renewable energy forecasting. *Energy Conversion and Management*, 198, 111799.
- [19] Lim, B., & Zohren, S. (2021). Time-series forecasting with deep learning: a survey. *Philosophical Transactions of the Royal Society A*, 379(2194), 20200209.
- [20] Torres, J. F., Hadjout, D., Sebaa, A., et al. (2021). Deep learning for time series forecasting: A survey. *Big Data*, 9(1), 3-21.
- [21] Wen, Q., Zhou, T., Zhang, C., et al. (2022). Transformers in time series: A survey. *arXiv preprint arXiv:2202.07125*.
- [22] Zhou, H., Zhang, S., Peng, J., et al. (2021). Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of AAAI* (Vol. 35, No. 12, pp. 11106-11115).
- [23] Wu, H., Xu, J., Wang, J., & Long, M. (2021). Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *Advances in Neural Information Processing Systems*, 34, 22419-22430.
- [24] Qin, Y., Song, D., Chen, H., et al. (2017). A dual-stage attention-based recurrent neural network for time series prediction. *arXiv preprint arXiv:1704.02971*.
- [25] Salinas, D., Flunkert, V., Gasthaus, J., & Januschowski, T. (2020). DeepAR: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, 36(3), 1181-1191.
- [26] Gasthaus, J., Benidis, K., Wang, Y., Rangapuram, S.S., Salinas, D., Flunkert, V., Januschowski, T. (2019). Probabilistic forecasting with spline quantile function RNNs. In *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics (AISTATS)*, PMLR 89: 1901-1910.

## 8. Annexes

### Annexe A : Spécifications techniques détaillées

#### A.1 Hardware de développement (entraînement)

- CPU : 13th Gen Intel(R) Core(TM) i5-13400F @ 4.0 GHz (10 cœurs, 16 threads)
- RAM : 64 GB DDR4-3200
- GPU : NVIDIA GeForce RTX 4070 (12 GB VRAM)
- Stockage : SSD NVMe KINGSTON SNV2S1000G (1 TB) + FIKWOT FX660 (2 TB)
- OS : Ubuntu 20.04 LTS
- CUDA : 11.2, cuDNN : 8.1
- TensorFlow : 2.8.0

#### A.2 Hardware de production (inférence edge)

- SBC : Raspberry Pi 4 Model B
- RAM : 8 GB LPDDR4-3200
- CPU : Broadcom BCM2711 quad-core Cortex-A72 @ 1.5 GHz
- TPU : Google Coral USB Accelerator
- Stockage : Carte SD SanDisk Extreme PRO 64 GB (UHS-I, U3, V30)
- Alimentation : 5V 3A USB-C officielle
- Refroidissement : Dissipateurs aluminium + ventilateur 5V
- OS : Raspberry Pi OS Lite 64-bit (Debian Bullseye)
- Python : 3.9.2
- TensorFlow Lite : 2.8.0
- PyCoral : 2.0.0

## Annexe B : Hyperparamètres optimisés

Configuration finale obtenue après grid search sur 128 configurations testées. La métrique d'optimisation est le MAPE sur ensemble de validation.

Tableau B.1 : Hyperparamètres du modèle LAD-BNet V7.1

Hyperparamètre	Valeur optimale	Alternatives testées
<b>Learning rate</b>	0.0005	0.0001, 0.001, 0.002
<b>Batch size</b>	16	8, 32, 64
<b>Dropout rate</b>	0.1	0.0, 0.2, 0.3
<b>Conv1D filters (1)</b>	64	32, 128
<b>Conv1D filters (2)</b>	64	32, 128
<b>Conv1D filters (dilated)</b>	128	64, 256
<b>Dense units (lag 1)</b>	256	128, 512
<b>Dense units (lag 2)</b>	128	64, 256
<b>Dense units (fusion 1)</b>	256	128, 512
<b>Dense units (fusion 2)</b>	128	64, 256
<b>Dilation rate</b>	2	1, 4, 8
<b>Kernel size</b>	3	5, 7
<b>Epochs</b>	400	Early stopping (patience 50)
<b>Optimizer</b>	Adam	RMSprop, SGD

Paramètres de quantification : Type int8 full integer, dataset représentatif de 1000 échantillons, quantization-aware training non utilisé (post-training quantization).

## Annexe C : Caractéristiques du dataset

Source : Système de monitoring énergétique du bâtiment IUT Réunion.

Tableau C.1 : Caractéristiques principales du dataset

Caractéristique	Valeur
Période	Janvier 2023 - Septembre 2024 (21 mois)
Résolution temporelle	10 minutes
Nombre total d'enregistrements	90,720
Variables brutes	datetime, DBT (°C), RH (%), kW
Taux de données manquantes	0.8%
Méthode d'imputation	Interpolation linéaire
Split entraînement	63,504 échantillons (70%)
Split validation	13,608 échantillons (15%)
Split test	13,608 échantillons (15%)

Statistiques descriptives des variables brutes :

Tableau C.2 : Statistiques descriptives

Variable	Minimum	Maximum	Moyenne	Écart-type	Médiane
DBT (°C)	15.2	32.8	24.3	3.7	24.1
RH (%)	32.0	95.0	68.5	12.3	70.0
kW	8.2	187.3	65.4	28.9	62.1

## Annexe D : Liste complète des features (28)

**Target (1):** kW : Consommation énergétique

**Météorologiques (2):** DBT : Température de bulbe sec (°C), RH : Humidité relative (%)

**Temporelles cycliques (6):** hour\_sin, hour\_cos, month\_sin, month\_cos, dayofweek\_sin, dayofweek\_cos

**Contextuelles (6):** weekend, is\_holiday, is\_business\_hours, is\_night, is\_morning\_peak, is\_evening\_peak

**Lags (5):** kW\_lag\_6, kW\_lag\_12, kW\_lag\_24, kW\_lag\_72, kW\_lag\_144

**Rolling statistics (6):** kW\_rolling\_mean\_6, kW\_rolling\_mean\_12, kW\_rolling\_mean\_24, kW\_rolling\_std\_12, kW\_rolling\_max\_24, kW\_rolling\_min\_24

**Interaction (1):** temp\_humidity\_interaction :  $DBT \times RH / 100$

**Total:** 28 features (27 en entrée + 1 target)

## Annexe E : Architecture du modèle (notation Keras)

```
# Architecture LAD-BNet en Python/Keras (simplifié)

from tensorflow.keras import Model, Input
from tensorflow.keras.layers import (Dense, Conv1D, Concatenate, Flatten,
                                      Lambda, BatchNormalization, Dropout,
                                      GlobalAveragePooling1D, GlobalMaxPooling1D)

# Input: (batch, 144 timesteps, 27 features)
inputs = Input(shape=(144, 27), name='input_sequence')

# === BRANCH 1: LAG BRANCH ===
lag_input = Lambda(lambda x: x[:, -24:, :])(inputs) # Derniers 24 pas
lag_flat = Flatten()(lag_input)
lag_d1 = Dense(256, activation='relu')(lag_flat)
lag_bn1 = BatchNormalization()(lag_d1)
lag_drop1 = Dropout(0.1)(lag_bn1)
lag_d2 = Dense(128, activation='relu')(lag_drop1)
lag_bn2 = BatchNormalization()(lag_d2)
lag_out = Dropout(0.1)(lag_bn2)

# === BRANCH 2: TCN BRANCH ===
tcn_c1 = Conv1D(64, 3, padding='causal', activation='relu')(inputs)
tcn_bn1 = BatchNormalization()(tcn_c1)
tcn_drop1 = Dropout(0.1)(tcn_bn1)
tcn_c2 = Conv1D(64, 3, padding='causal', activation='relu')(tcn_drop1)
tcn_bn2 = BatchNormalization()(tcn_c2)
tcn_drop2 = Dropout(0.1)(tcn_bn2)
tcn_dilated = Conv1D(128, 3, padding='causal', dilation_rate=2,
                     activation='relu')(tcn_drop2)
tcn_bn3 = BatchNormalization()(tcn_dilated)
tcn_drop3 = Dropout(0.1)(tcn_bn3)
tcn_avg = GlobalAveragePooling1D()(tcn_drop3)
tcn_max = GlobalMaxPooling1D()(tcn_drop3)
tcn_out = Concatenate()([tcn_avg, tcn_max])

# === FUSION MODULE ===
merged = Concatenate()([lag_out, tcn_out])
fus_d1 = Dense(256, activation='relu')(merged)
fus_bn1 = BatchNormalization()(fus_d1)
fus_drop1 = Dropout(0.1)(fus_bn1)
fus_d2 = Dense(128, activation='relu')(fus_drop1)
fus_drop2 = Dropout(0.1)(fus_d2)
outputs = Dense(72, name='output_predictions')(fus_drop2)

# Modèle final
model = Model(inputs=inputs, outputs=outputs, name='LAD_BNet')
model.compile(optimizer='adam', loss='mse', metrics=['mae', 'mape'])

# Paramètres totaux: ~245,000
# Taille float32: ~12 MB
# Taille int8 quantifié: ~3 MB
```

## **Annexe F : Disponibilité et licence / Availability and License**

### **Licence de l'article / Article License**

**Cet article est publié sous licence Creative Commons Attribution 4.0 International (CC-BY 4.0).**

**This article is licensed under Creative Commons Attribution 4.0 International (CC-BY 4.0). <https://creativecommons.org/licenses/by/4.0/>**

**Cette licence vous permet de / You are free to:**

- Partager — copier, distribuer et communiquer le matériel**  
**Share — copy and redistribute the material in any medium or format**
- Adapter — remixier, transformer et créer à partir du matériel**  
**Adapt — remix, transform, and build upon the material for any purpose**

**Sous réserve de / Under the following terms:**

- Attribution — Vous devez créditer l'œuvre et indiquer si des modifications ont été effectuées**

**Attribution — You must give appropriate credit and indicate if changes were made**

### **Licence du code source / Source Code License**

**Code source / Source code: MIT License**

**Disponibilité / Availability:**

- Code source : Disponible sur demande (dépôt GitHub en cours de préparation)**  
**Source code: Available upon request (GitHub repository in preparation)**
- Modèles pré-entraînés : Disponibles sur demande pour recherche académique**  
**Pre-trained models: Available upon request for academic research**
- Datasets : Données anonymisées disponibles sur demande (respectant RGPD)**  
**Datasets: Anonymized data available upon request (GDPR compliant)**
- Documentation technique complète : Incluse avec le code source**  
**Complete technical documentation: Included with source code**

**Contact : [jean-philippe.lignier@tangibleassets.org](mailto:jean-philippe.lignier@tangibleassets.org)**

**Licence : MIT License**