

A SYNTHETIC DATASET OF FRENCH ELECTRIC LOAD CURVES WITH TEMPERATURE CONDITIONING

Tahar Nabil¹, Ghislain Agoua¹, Pierre Cauchois², Anne De Moliner², Benoît Grossin¹

¹ EDF R&D, Palaiseau, France

² Enedis, France

{tahar.nabil, ghislain.agoua, benoit.grossin}@edf.fr

ABSTRACT

The undergoing energy transition is causing behavioral changes in electricity use, e.g. with self-consumption of local generation, or flexibility services for demand control. To better understand these changes and the challenges they induce, accessing individual smart meter data is crucial. Yet this is personal data under the European GDPR. A widespread use of such data requires thus to create synthetic realistic and privacy-preserving samples. This paper introduces a new synthetic load curve dataset generated by conditional latent diffusion. We also provide the contracted power, time-of-use plan and local temperature used for generation. Fidelity, utility and privacy of the dataset are thoroughly evaluated, demonstrating its good quality and thereby supporting its interest for energy modeling applications.

1 INTRODUCTION

Meeting the Net Zero Emissions by 2050 Scenario calls for a swift energy transition, involving more renewables on the production side and an intense electrification on the demand side (IEA, 2024). Optimizing supply and balance, developing new energy services; e.g. around self-consumption of local renewable production or demand control, become thus both increasingly challenging and dependent on one key resource: access to fine-grained electricity consumption data at the individual level (Chen et al., 2023). However, strict privacy regulations severely limit the sharing of smart meter data (Enedis, 2024). This creates a significant barrier for the research community to develop innovative models for load curve analysis.

On the other hand, the remarkable successes of deep learning generative models in several domains has contributed to initiate a promising area of research: learning to generate synthetic yet realistic smart meter data, while ensuring their privacy (Chai et al., 2024). Nevertheless, the available synthetic datasets suffer from certain shortcomings. Specifically, they are either (i) coarse-grained at the hourly timestep and building level (Emami et al., 2023), lacking the richness observed at the individual sub-hourly levels; or (ii) restricted to daily profile data (Liang and Wang, 2022; Yuan et al., 2023; Thorve et al., 2023; Chai and Chadney, 2024) with no conditioning on temperature, whereas the latter has a significant impact on demand in countries where electric space heating (e.g. France) or cooling prevails.

In this work, we release a new synthetic dataset of one-year 30-min resolution electric load curves by leveraging latent diffusion models (Rombach et al., 2022). We extend their conditioning mechanism to time-varying exogenous variables. This allows us to condition our samples on outdoor temperature, also released. A thorough evaluation shows that the proposed model (1) yields unprecedented high-quality samples, with superior performance compared to TimeGAN, a standard time series generative model, without compromising privacy; (2) can flexibly incorporate any static information describing the consumer, enabling us to feature a wide range of customer contracted powers and individual behaviors.

2 CONDITIONAL LOAD CURVE GENERATION WITH LATENT DIFFUSION

Our dataset is generated by training a conditional latent diffusion model (Rombach et al., 2022), which we describe next. A brief discussion of related public smart meter datasets is provided in Appendix A.1. Our architecture is depicted in Supplementary Figure 2.

Latent diffusion model (LDM) Load curves $\mathbf{x} \in \mathbb{R}^{T \times 1}$, with $T = N_{days} \times 48$ at sampling rate 30 minutes, are seen as single-channel 2D images $\tilde{\mathbf{x}} \in \mathbb{R}^{1 \times N_{days} \times 48}$. LDMs are two-stage models. (1) First, we fit a convolutional autoencoder with compression factor $r = 4$ to reconstruct $\tilde{\mathbf{x}}$. \mathbf{x} is therefore represented by a latent code $\mathbf{z} \in \mathbb{R}^{c_z \times N_{days}/r \times 48/r}$ of c_z channels. The loss function of the autoencoder adds a vector quantization term to the L2 reconstruction loss to regularize the latent space (Rombach et al., 2022; van den Oord et al., 2017). (2) Next, the distribution $p(\mathbf{z})$ of codes in the latent space of the frozen autoencoder is approximated by learning a Denoising Diffusion Probabilistic Model (DDPM, Ho et al. (2020)). Specifically, a UNet with spatial attention is trained to minimize the denoising objective.

Conditioning mechanisms Leveraging the flexibility of latent diffusion, we propose to condition the model on both static variables and other time-varying exogenous variables.

1. *Static variables*: following Rombach et al. (2022), the UNet model is conditioned on labels by either concatenation to \mathbf{z} or cross-attention through a spatial transformer.
2. *Exogenous variables*: to condition \mathbf{x} on another series $\mathbf{u} \in \mathbb{R}^{T \times 1}$, we process its code \mathbf{z} with a cross-attention layer, with \mathbf{z} as the query and patched $\mathbf{u} \in \mathbb{R}^{(T//P) \times P}$ as the keys and values, where P is the patch length (Nie et al., 2023). This results in a transformed code $\mathbf{h} \in \mathbb{R}^{c_z \times N_{days}/r \times 48/r}$ which is then processed by the decoder.

Hence (i) conditioning on exogenous variables such as outdoor temperature is handled by the decoder only and does not affect the diffusion loss for learning the distribution of load curves; but (ii) conditioning on static variables does not affect fitting the autoencoder for obtaining good reconstructions. Provided that the autoencoder generalizes correctly, any static conditioning can thus be added in the second stage of learning the data distribution.

Generation At inference, the diffusion process is reversed to sample a code $\hat{\mathbf{z}}$ conditionally on labels. $\hat{\mathbf{z}}$ undergoes cross-attention with the exogenous variables and is then decoded to data space.

3 EVALUATION

We compare our latent diffusion model (LDM) to TimeGAN, the standard baseline for time series generation (Yoon et al., 2019). Both are trained on a set of 17k one-year French load curves, conditionally on temperature and on two labels of 3 classes each: contracted power and time-of-use (ToU) rate. They are evaluated along the three key dimensions of synthetic data: fidelity (Section 3.1), utility (in 3.2) and privacy (in 3.3). Full data description, models implementation, metrics and results are given respectively in Appendices A.2, B, C and D.

3.1 FIDELITY & DIVERSITY

The evaluation relies on three metrics. (i) The discriminative score (Yoon et al., 2019) is the test accuracy of a *post-hoc* 1-NN binary classifier separating real and synthetic data with Euclidean distance. (ii) Context-FID measures how well the synthetic dataset recovers the training set statistics in a suitable embedding space (Jeha et al., 2022; Franceschi et al., 2019). (iii) The correlation score measures temporal consistency (Ni et al., 2022).

Results The metrics shown in Table 1 and Supplementary Table 3, demonstrate that LDM consistently outperforms TimeGAN, with high-quality samples challenging to distinguish from real data. Qualitatively, (i) the 2D-projection of the samples with t-SNE in Figure 1(a) is consistent with the discriminative scores, with LDM samples hard to distinguish from real samples; (ii) LDM outperforms TimeGAN as shown by the histograms of daily statistics of interest (mean, max, quantiles,

Table 1: Fidelity scores on the hold-out test set for samples on *night* ToU and 6 kVA. Discriminative score computed on one-year load curves (D_{year}) or on averaged daily profiles ($D_{profile}$). Best performance emphasized in bold. Full results in Supplementary Table 3.

	D_{year} (\downarrow)	$D_{profile}$ (\downarrow)	Context-FID (\downarrow)	Correlation score (\downarrow)
LDM	0.037	0.059	1.748	0.002
TimeGAN	0.357	0.452	2.082	0.224

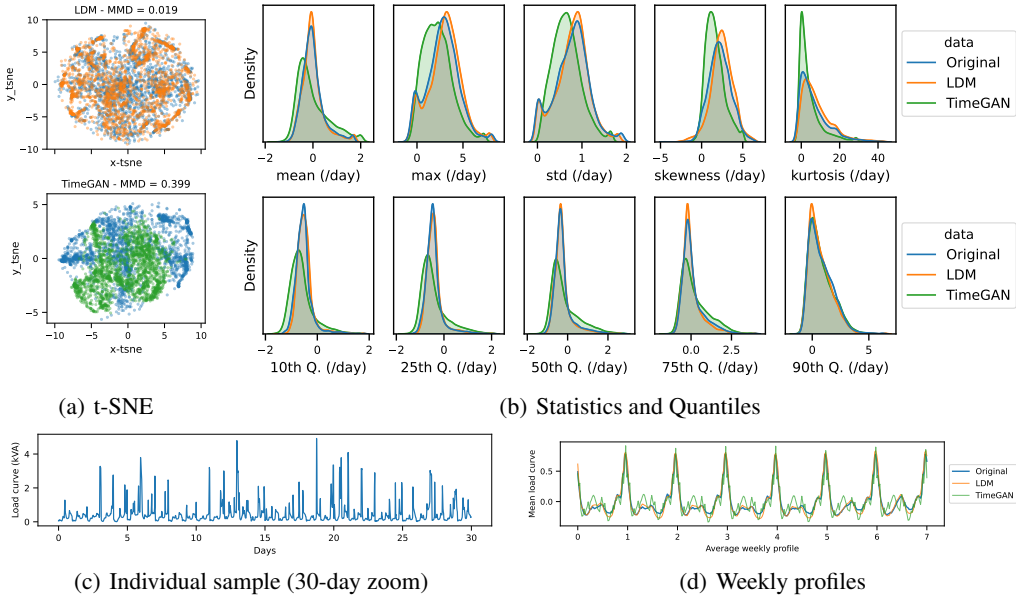


Figure 1: (a) t-SNE 2D projection of original (blue) and synthetic data across all categories for Latent Diffusion (orange) and TimeGAN (green). (b) to (d) are restricted to the night ToU, 6kVA category: (b) Density estimation of daily statistics and quantiles; (c) Example of a synthetic sample by Latent Diffusion; (d) Mean weekly profiles.

etc.) in Figure 1(b), by the average weekly profiles in Figure 1(d), as well as by additional plots in Appendix D.1.1. (iii) Some LDM individual samples are displayed in Figures 1(c) and 8-10, showcasing its ability to generate diverse patterns, e.g. plateaus of low consumption, while respecting known seasonalities (daily peaks).

Thermo-sensitivity We further evaluate whether the synthetic data preserve the correlation between temperature and consumption. Quantitatively, we compare the distributions of the thermo-sensitivity gradients, which explain daily variations in consumption in winter by those in temperature (see Appendix A.3). Qualitatively, we generate load curves under a strong global temperature offset throughout the year. The results, in Appendix D.1.2, suggest that the latent diffusion model learned a meaningful conditioning by temperature, with (i) gradients matching the real test distribution and (ii) realistic distortion of the load curves under the temperature offset.

3.2 UTILITY

Utility assesses whether generated data can be used to solve machine learning tasks to the same extent as with real data. Beyond next-step prediction (Yoon et al., 2019), new tasks must be designed for electricity consumption (Chai et al., 2024). We adopt the standard *Train on Synthetic, Test on Real* (TSTR) setting, where *post-hoc* models are trained either on synthetic data or real train data, and compared on a hold-out real test set. We tackle two tasks, short-term load forecasting and time series classification.

Table 2: TSTR metrics; the closest model to TRTR (*Train on Real, Test on Real*) is emphasized in bold. Forecasting results are averaged across horizons [48, 96, 192, 336], for a lookback of length 720. Baselines: copy from last week (forecasting), majority class (classification).

Loss	TRTR	LDM	TimeGAN	Baseline
Forecast. MSE	0.190	0.190	0.209	0.306
Forecast. MAE	0.234	0.233	0.253	0.251
Classif. Acc	0.740	0.750	0.700	0.607
Classif. F1	0.576	0.564	0.537	0.252

Forecasting Using a lookback-window of 15 days, we train a *post-hoc* state-of-the-art transformer model, PatchTST (Nie et al., 2023), to predict individual load curves for various horizons H , from a day ($H = 48$) to a week ($H = 336$). Table 2 shows the mean squared (MSE) and absolute (MAE) errors, averaged across all horizons (detailed scores in Supplementary Table 4): models trained on either real or latent diffusion data are in close agreement, whereas training on TimeGAN degrades the MSE on real test data by about 9%.

Classification We train a kNN ($k=5$) classifier to predict the ToU (3 classes) associated to each load curve. The input representation of the classifier is a 104-dimensional vector, where each load curve is transformed into a vector of daily profile over the year (48 values), daily profile restricted to winter (48 values) and a vector of 8 statistics. The ground truth for fake data is the conditioning label. The test accuracy and F1-scores in Table 2 yield conclusions similar to the forecasting task, confirming that LDM produces high-quality data.

3.3 PRIVACY

We evaluate the privacy of synthetic samples through two Membership Inference Attacks (MIA, Carlini et al. (2022)) and a statistical test of relative similarity (MMD-test).

MIA Commonly considered as the first step towards auditing the privacy of deployed models, MIAs try to predict whether a given sample has been seen by the target model during its training. We perform two attacks, described in Appendix C: (i) black-box, by exploiting the distance to generated samples; (ii) white-box, by exploiting the reconstruction scores of the autoencoders. Following best practices (Carlini et al., 2022), we report the ROC curve in log-scale in Supplementary Figure 11 and the true positive rate at 0.1% false positive rate in Supplementary Table 5. The performance of the attacks are close to random for both models, at every false positive rate.

MMD-test In addition, we perform a three-sample MMD test with the following null hypothesis: the distribution of synthetic daily profiles is closer to test than train data (Bounliphone et al., 2016). The p-values in Supplementary Table 6 are large (> 0.2 for all categories but one at 0.019), suggesting no evidence for the models to overfit training data. Similar conclusions apply to TimeGAN.

Together with the computation of the nearest neighbor distance ratio (Appendix D.3), these tests converge thus to the same conclusion: without being a formal proof of privacy, they are hints that the model is not prone to overfitting and yields original synthetic samples, not copying the train set.

4 CONCLUSION AND IMPACT STATEMENT

This work proposes a new and richer synthetic dataset for residential load curves, covering one year, with contracted power and ToU plan information, as well as the local temperature data used for generation. The generation is achieved by leveraging the flexibility of latent diffusion models and an astute conditioning mechanism that incorporates both static and dynamic conditional variables. The evaluation shows good fidelity, utility and privacy metrics. Most notably, LDM consistently achieves high-quality scores in every fidelity aspect.

Future work could extend the generation to more categories, e.g. non-thermo-sensitive load curves or small and medium-sized businesses. Besides, properly assessing the utility of smart meter data

remains an open question: other ideas could exploit the conditioning on temperature to forecast with exogenous variables (Wang et al., 2024), perform transfer experiments (Emami et al., 2023) or train time series foundation models (Woo et al., 2024).

DATA AVAILABILITY

We publicly release a synthetic dataset (doi: [10.5281/zenodo.15232742](https://doi.org/10.5281/zenodo.15232742)) of 10k one-year residential load curves, with their local temperatures, contracted power and ToU plans. This open dataset generated by conditional latent diffusion captures the correlation between cold temperature and electric consumption that can be found in certain countries, e.g. in France.

ACKNOWLEDGEMENT

This paper is based on joint work between Enedis and EDF R&D in accordance with regulated R&D contract. Moreover, we would like to thank Etienne Le Naour for the valuable discussions on this project and feedbacks on the paper.

REFERENCES

- T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2019.
- W. Bounliphone, E. Belilovsky, M. B. Blaschko, I. Antonoglou, and A. Gretton. A test of relative similarity for model selection in generative models. In *The Fourth International Conference on Learning Representations*, 2016.
- N. Carlini, S. Chien, M. Nasr, S. Song, A. Terzis, and F. Tramèr. Membership inference attacks from first principles. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 1897–1914, 2022. doi: 10.1109/SP46214.2022.9833649.
- S. Chai and G. Chadney. Faraday: Synthetic smart meter generator for the smart grid. In *ICLR 2024 Workshop on Tackling Climate Change with Machine Learning*, 2024. URL <https://www.climatechange.ai/papers/iclr2024/43>.
- S. Chai, G. Chadney, C. Avery, P. Grunewald, P. Van Hentenryck, and P. L. Donti. Defining 'good': Evaluation framework for synthetic smart meter data. *arXiv preprint arXiv:2407.11785*, 2024.
- Z. Chen, A. M. Amani, X. Yu, and M. Jalili. Control and optimisation of power grids using smart meter data: A review. *Sensors*, 23(4), 2023. ISSN 1424-8220. doi: 10.3390/s23042118.
- K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, 2014. doi: 10.3115/v1/D14-1179.
- P. Emami, A. Sahu, and P. Graf. Buildingsbench: A large-scale dataset of 900k buildings and benchmark for short-term load forecasting. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=c5rqd6PZn6>.
- Enedis. Données personnelles, 2024. <https://www.enedis.fr/donnees-personnelles> [In French, accessed on January 14, 2025].
- J.-Y. Franceschi, A. Dieuleveut, and M. Jaggi. Unsupervised Scalable Representation Learning for Multivariate Time Series. In *Advances in Neural Information Processing Systems*, volume 33, pages 4650–4661, 2019. URL <https://proceedings.neurips.cc/paper/2019/file/53c6de78244e9f528eb3e1cda69699bb-Paper.pdf>.
- J. Hayes, L. Melis, G. Danezis, and E. De Cristofaro. LOGAN: Membership inference attacks against generative models. *arXiv preprint arXiv:1705.07663*, 2017.

- J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, pages 6840–6851. Curran Associates, Inc., 2020.
- IEA. Global energy and climate model, iea, paris, 2024. <https://www.iea.org/reports/global-energy-and-climate-model> [Accessed on January 17, 2025].
- P. Jeha, M. Bohlke-Schneider, P. Mercado, S. Kapoor, R. S. Nirwan, V. Flunkert, J. Gasthaus, and T. Januschowski. PSA-GAN: Progressive self attention GANs for synthetic time series. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=Ix_mh42xq5w.
- T. Kim, J. Kim, Y. Tae, C. Park, J.-H. Choi, and J. Choo. Reversible instance normalization for accurate time-series forecasting against distribution shift. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=cGDAkQo1C0p>.
- X. Liang and H. Wang. Synthesis of realistic load data: Adversarial networks for learning and generating residential load patterns. In *NeurIPS 2022 Workshop on Tackling Climate Change with Machine Learning*, 2022. URL <https://www.climatechange.ai/papers/neurips2022/93>.
- I. Loshchilov and F. Hutter. Decoupled Weight Decay Regularization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- H. Ni, L. Szpruch, M. Sabate-Vidales, B. Xiao, M. Wiese, and S. Liao. Sig-Wasserstein GANs for time series generation. In *Proceedings of the Second ACM International Conference on AI in Finance*, ICAIF '21, 2022. doi: 10.1145/3490354.3494393.
- Y. Nie, N. H. Nguyen, P. Sinthong, and J. Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=Jbdc0vTOcol>.
- R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- S. Thorve, Y. Y. Baek, S. Swarup, H. Mortveit, A. Marathe, A. Vullikanti, and M. Marathe. High resolution synthetic residential energy use profiles for the United States. *Scientific Data*, 10(1): 76, 2023. ISSN 2052-4463. doi: 10.1038/s41597-022-01914-1.
- A. Trindade. ElectricityLoadDiagrams20112014. UCI Machine Learning Repository, 2015. DOI: <https://doi.org/10.24432/C58C86>.
- A. van den Oord, O. Vinyals, and k. kavukcuoglu. Neural discrete representation learning. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- Y. Wang, H. Wu, J. Dong, G. Qin, H. Zhang, Y. Liu, Y.-Z. Qiu, J. Wang, and M. Long. Timexer: Empowering transformers for time series forecasting with exogenous variables. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=INAEUQ041T>.
- G. Woo, C. Liu, A. Kumar, C. Xiong, S. Savarese, and D. Sahoo. Unified training of universal time series forecasting transformers. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=Yd8eHMY1wz>.
- J. Yoon, D. Jarrett, and M. Van der Schaar. Time-series generative adversarial networks. In *Advances in Neural Information Processing Systems*, volume 32, 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/hash/c9efe5f26cd17ba6216bbe2a7d26d490-Abstract.html.
- R. Yuan, S. A. Pourmousavi, W. L. Soong, A. J. Black, J. A. R. Liisberg, and J. Lemos-Vinasco. A synthetic dataset of Danish residential electricity prosumers. *Scientific Data*, 10(1):371, 2023. ISSN 2052-4463. doi: 10.1038/s41597-023-02271-3.

A BACKGROUND

A.1 RELATED WORK

Open electricity consumption data First open electricity consumption (load curve) datasets typically contained real smart meter readings for the aggregated household consumption at sub-hourly timesteps, e.g. 15 minutes (Electricity, in Portugal, [Trindade \(2015\)](#)) or 30 minutes ([Irish Commission for Energy Regulation dataset](#), [UK Power Networks \(London\)](#), [Smart-Grid Smart-City Customer Trial Data, Australia](#)). However, they are of limited use to grid operators and others involved in the energy transition, dating back to around a decade ago and corresponding thus to other user behaviors. Similarly, [conso-inf36](#) is an open dataset of regularly updated 30-min resolution French load curves, but limited to averages of at least 100 individuals. Recent works have investigated the use of generative models to produce synthetic realistic consumption data, e.g. ([Liang and Wang, 2022](#); [Yuan et al., 2023](#); [Thorve et al., 2023](#); [Chai and Chadney, 2024](#)). These references restrict themselves to daily profiles at hourly or sub-hourly timesteps. For instance, Faraday generates daily profiles with a Variational Auto-Encoder and provides information such as property type and low carbon technologies ownership (e.g. electric vehicles ownership) ([Chai and Chadney, 2024](#)). A common limitation of all aforementioned references is the lack of conditioning by outdoor temperature. Finally, another line of work relies on simulation: BuildingsBench ([Emami et al., 2023](#)) simulates a large number of one-year hourly residential and commercial building load curves, representative of the US stock, conditionally on local temperature – without releasing the temperature time series.

A.2 DATA DETAILS

Training data The available training dataset contains the following variables:

- 17k residential load curves spanning 94 departments in metropolitan France. The individuals have time-of-use (ToU) plan, their consumption is correlated with outdoor temperature, particularly during winter. Data is at 30-min resolution and covers one year starting from October 2022;
- Local outdoor temperature data for the same period;
- Contracted power information: 6, 9, 12 (kVA);
- Time-of-use (ToU) plan: 3 classes "midday", "night" and "misc" describe the time of day when prices are lower. ToU strongly conditions the load curve, in particular the daily patterns and the time of peak consumption.

Test data contain 2k samples with the same information as for training.

A.3 THERMO-SENSITIVITY GRADIENTS

When aggregated at a daily granularity, the analysis of load data and temperature shows an almost linear relationship between load and temperature below a temperature threshold. The thermo-sensitivity gradient represents the relationship between load and temperature below the threshold. The steps to calculate the gradient for a daily load curve are as follows:

1. filter on winter data (from November to March) to avoid bias;
2. calculate the load delta: $\Delta Load_d = Load_d - Load_{d-7}$ where $Load_d$ is the load for day d in kWh;
3. calculate the degree-day $DJU_d = \max(0, T_{thresh} - T_d)$ where T_d is the temperature for day d and $T_{thresh} \in [14.5, 18]$ (in °C) the temperature threshold for the region;
4. calculate the delta of degree-days $\Delta DJU_d = DJU_d - DJU_{d-7}$;
5. compute the linear regression $\Delta Load = Gradient \times \Delta DJU$.

The gradients are in kWh/degree-day.

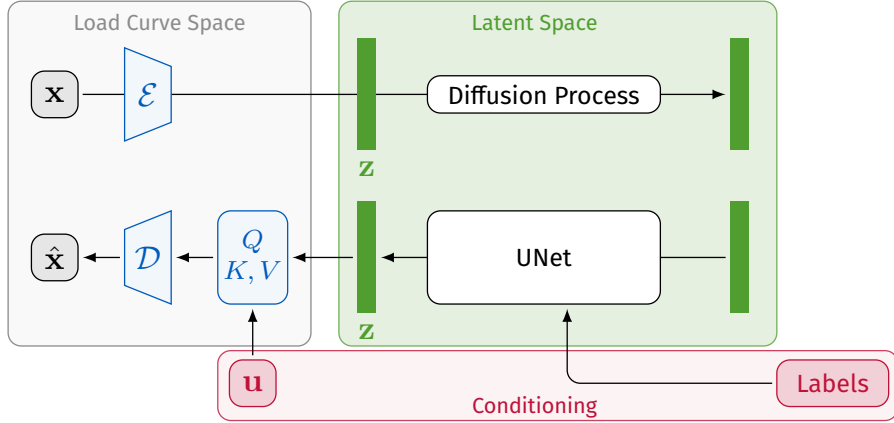


Figure 2: Latent diffusion (Rombach et al., 2022) with conditioning on exogenous variables. \mathcal{E} : image encoder, \mathcal{D} : image decoder, Q, K, V : cross-attention network between latent vector \mathbf{z} and exogenous time series \mathbf{u} . \mathbf{x} is the input data (load curve) and $\hat{\mathbf{x}}$ its reconstruction. Static labels are added by cross-attention to the diffusion UNet, or by concatenation to the latent codes \mathbf{z} .

B MODEL AND BASELINE DETAILS

B.1 LATENT DIFFUSION

Architecture and hyperparameters The overall architecture is represented in Figure 2. The hyperparameters of the first stage autoencoder model are:

- Compression ratio: $r = 4$;
- Input channels: 1;
- Number of channels: 128, channel multipliers: [1, 2, 4];
- Codebook: number of channels $c_z = 3$, shape $3 \times 90 \times 12$, cardinality: 8192;
- Temperature conditioning:
 - Temperature (in $^{\circ}\text{C}$) is scaled by 35°C and reshaped as non-overlapping patches of size 32;
 - Cross-attention network: attention heads: 8, attention layers: 4, positional encoding with sines and cosines;
- Model size: 61M parameters.

The hyperparameters for the denoising diffusion model are:

- Diffusion steps: 1000;
- Noise schedule: linear;
- UNet architecture (spatial attention):
 - Attention heads: 16;
 - Number of channels: 160, multipliers: [1, 2, 3, 4];
- Label representation: contracted power normalized between 0 and 1 is concatenated to a learned embedding of dimension 4 of the tariff (3 classes);
- Label conditioning: concatenation to the latent code along the channel dimension;
- Model size: 157M parameters.

Implementation details Load curves are z-normalized instance-wise, by computing the mean and standard deviation across the time dimension, for each sample. The autoencoder is trained with Huber loss, a learning rate of 5×10^{-5} for 250k optimization steps of batch size 16 and a cosine annealing schedule with 10k warmup steps. We also add data augmentation when training the autoencoder: with probability 0.5, we shift the temperature by a random global offset δ , and modify consequently the load curve by adding $-g \times \delta$, with $g > 0$ a random thermo-sensitivity gradient. The diffusion model is trained with mean squared error loss for the denoising objective (see [Ho et al. \(2020\)](#); [Rombach et al. \(2022\)](#)), for 200k steps with batch size 16, learning rate 3×10^{-6} , cosine annealing with 10k warmup steps. For both the autoencoder and the diffusion model, we used the AdamW algorithm to optimize the loss functions ([Loshchilov and Hutter, 2019](#)), with default parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$, $\lambda = 0.01$. All models were trained on a single NVIDIA A100-40GB GPU.

B.2 TIMEGAN

The TimeGAN model ([Yoon et al., 2019](#)) is a Generative Adversarial Network model that operates in the latent space of an autoencoder. Therefore, it consists of 4 networks: encoder, decoder, generator and discriminator, which are all GRUs ([Cho et al., 2014](#)) to better handle the temporal dynamics. Similarly to Latent Diffusion, each time series $\mathbf{x} \in \mathbb{R}^{T \times 1}$ is represented as a sequence of N_{days} 48-dimensional vectors $\tilde{\mathbf{x}} \in \mathbb{R}^{N_{days} \times 48}$, with $T = N_{days} \times 48$. Conditioning by temperature is achieved by concatenation of the temperature to the inputs of all 4 sub-networks of TimeGAN, along the time dimension. Similarly for conditioning by discrete labels, which are repeated N_{days} times and stacked along the second time dimension as well. The input load curve data is transformed through an inverse hyperbolic sine (*asinh*) transformation. Load curve and temperature are then normalized with a min-max scaler. The parameters of the models were optimized using Optuna ([Akiba et al., 2019](#)) and grid search. The final parameters are:

- hidden dimension $h_{dim} = 32$ for the latent space of the autoencoder, $z_{dim} = 24$ for the generator and $n_{layers} = 1$ layer in all GRU's networks;
- we use different learning rates $lr = 0.0005$ for the encoder, $lr_{gen} = 0.0005$ for the generator, and $lr_{disc} = 0.005$ for the discriminator to increase the stability of the training;
- the model is trained with a batch size of 32 and 15k steps.

C EVALUATION PROTOCOL

Description The evaluation of the generative models relies on a standard protocol, by comparing sets of synthetic data to a hold-out test set of actual load curves. We generate as many samples per conditioning label as there is in the test set, and use the temperatures in the test set: both synthetic and test sets have the same number of samples n_{test} . The two sets are compared along the three dimensions: fidelity, utility and privacy.

Fidelity metrics Fidelity metrics assess whether the synthetic data are realistic and diverse. The following metrics are computed:

- *Discriminative score*: following [Yoon et al. \(2019\)](#), we use $n_{test}/2$ original and $n_{test}/2$ synthetic samples to train a binary classifier, where original (respectively, synthetic) data are labelled 0 (resp., 1). We compute the accuracy acc on the rest of the original and generated data and report $|1/2 - acc|$ as the discriminative score. Perfectly separable original and synthetic sets are expected to achieve a score of 0.5; whereas indistinguishable sets are close to 0. We use a 1-nearest-neighbor with Euclidean distance as classifier, computed either in the space of yearly data (D_{year}) or after averaging the load curves to daily profiles of 48 values ($D_{profile}$).
- *Context-FID* is a Frchet-Inception distance-like score for time series. In computer vision, FID assesses the quality of images by comparing the distribution of synthetic vs real images, through the means and covariance matrices of their respective embeddings, produced by an Inception network. Here, following [Jeha et al. \(2022\)](#), the embedding network consists in temporal convolutions with triplet loss (contrastive learning with an encoder-only

architecture, from [Franceschi et al. \(2019\)](#)). We train one network per dataset, synthetic or real. To better embed local contexts, we split the one-year series into one-month series and fit one model per month. We report the average of the FID scores computed for all 12 months. Best scores are the lowest.

- The hyperparameters for training the encoder network are: 10k steps with batch size 128, learning rate = 10^{-4} , 10 layers of dilated causal 1D convolutions with 40 channels, output convolution with 320 channels, linear projection to a latent dimension of 160. The final max-pooling layer across sequence dimension guarantees that all series are encoded as vectors in \mathbb{R}^{160} .
- *Correlation score*: assesses temporal consistency by computing the absolute difference between auto-correlation functions of real and synthetic data ([Ni et al., 2022](#)).

Utility metrics Utility metrics measure whether synthetic data are useful for solving other downstream machine learning tasks. They are computed in the *TSTR*, *Train on Synthetic, Test on Real* scenario, i.e by training a model on synthetic data and testing it on a hold-out test set of real data. The test metrics are compared to the same model trained on real data (*TRTR*, *Train on Real, Test on Real*).

- Forecasting task: the hyperparameters of the forecasting model, PatchTST ([Nie et al., 2023](#)) are as follows:
 - Lookback window: 720, horizon $H \in [48, 96, 192, 336]$
 - Patch length: 32, stride: 16
 - Attention: 3 layers, 16 heads and $d_{model} = 128$ (token dimension)
 - Model trained with reversible instance normalization (RevIN, [Kim et al. \(2022\)](#))
 - Learning rate: 10^{-6}
 - Batch size: 256, number of steps: 20k, cosine learning rate scheduler with 1k warmup steps

For fairness, all training sets in TSTR or TRTR are restricted to the same size. TRTR results are averaged over 5 runs from random subsamples of the training set.

Privacy metrics

- *Black-box membership inference attack*: Membership inference attacks are fundamental privacy attacks of machine learning models, trying to predict if a particular sample has been seen during training. They are the foundations for stronger attacks ([Carlini et al., 2022](#)). In a black box scenario, the attacker has only access to queries to the target model under attack, and is given a set $\mathbf{x}_1, \dots, \mathbf{x}_{n+m}$ of $n + m$ samples of which n are known to be from the training set. We perform a simple attack: the attacker decides to classify as train the samples with score below a threshold, and as test otherwise. In this blackbox scenario, the score is the minimum distance between the target sample and any synthetic sample.
- *White-box (loss) membership inference attack*: assumes that the attacker has access to the trained machine learning model - although in our case, the model is not publicly available. The attack is similar to blackbox, except that the score is the reconstruction error of the samples by the autoencoder ([Hayes et al., 2017](#)).
- *Three-sample Maximum Mean Discrepancy Test*: initially, [Bounliphone et al. \(2016\)](#) introduced a statistical test of relative similarity to determine which of two models generates samples that are significantly closer to a real-world reference dataset of interest. The test statistics are based on differences in maximum mean discrepancies (MMDs). Here, we use the test to assess whether the train set is significantly closer to the generated samples than the test set. In other words, our three samples are: the synthetic dataset and the original train and test sets. We perform the test in the space of daily profiles instead of one-year curves because of the cost of MMD computation and since profiles are easier to discriminate. The null hypothesis is thus \mathcal{H}_0 : the distribution of synthetic daily profiles is closer to test daily profiles than train daily profiles. See ([Bounliphone et al., 2016](#)) for more details on the implementation of the test. Rejecting \mathcal{H}_0 would imply that there is statistical

Table 3: Fidelity scores on the hold-out test set for all categories of contracted power and time-of-use (ToU). Discriminative score computed on the entire load curves (D_{year}) or on averaged daily profiles ($D_{profile}$). Best performance emphasized in bold.

Contracted Power	ToU	Model	D_{year} (\downarrow)	$D_{profile}$ (\downarrow)	Context-FID (\downarrow)	Correlation score (\downarrow)
6 kVA	midday	LDM	0.049	0.078	2.394	0.003
		TimeGAN	0.388	0.471	2.802	0.130
	night	LDM	0.037	0.059	1.748	0.002
		TimeGAN	0.357	0.452	2.082	0.224
	misc.	LDM	0.119	0.0	3.405	0.013
		TimeGAN	0.357	0.405	4.021	0.113
9 kVA	midday	LDM	0.035	0.054	2.000	0.059
		TimeGAN	0.381	0.438	2.483	0.021
	night	LDM	0.068	0.051	1.622	0.022
		TimeGAN	0.414	0.434	1.968	0.127
	misc.	LDM	0.030	0.091	4.219	0.022
		TimeGAN	0.424	0.379	4.631	0.029
12 kVA	midday	LDM	0.024	0.032	2.913	0.011
		TimeGAN	0.395	0.379	3.301	0.048
	night	LDM	0.071	0.150	1.968	0.012
		TimeGAN	0.387	0.444	2.278	0.089
	misc.	LDM	0.0	0.043	4.744	0.005
		TimeGAN	0.370	0.370	5.075	0.038
all	all	LDM	0.029	0.026	1.235	0.011
		TimeGAN	0.382	0.432	1.476	0.123

evidence that the generative model has overfitted the train set, which would be a breach of privacy.

- *Nearest Neighbor Distance Ratio (NNDR)*: given a real train (respectively, test) sample, we compute the NNDR as the ratio of (i) the distance to its nearest neighbor in the synthetic set to (ii) the distance to its second nearest neighbor, in the synthetic set. Values of NNDR close to 1 suggest that synthetic samples are located in dense regions of real data, whereas values close to 0 indicate that some synthetic samples are located in sparse regions, i.e. near outliers.

D FULL EXPERIMENTAL RESULTS

D.1 FIDELITY

D.1.1 METRICS AND PLOTS

Metrics We provide the breakdown of fidelity metrics for all 9 categories of ToU \times contracted power, as well as for the entire synthetic dataset, in Table 3. Latent Diffusion consistently achieves high-quality discriminative scores in all categories (close to 0 in all categories, and less than 0.07 in 9/10 cases). The gap is significant with TimeGAN, which fails to produce samples undistinguishable from the real distribution. The Context-FID and correlation scores are also strongly in favor of Latent Diffusion over TimeGAN, supporting the former as a better generative model.

Plots This section provides additional plots at the aggregated level, i.e. by computing averages over groups of samples. First, the average one-year load curves is shown in Figure 3, for one conditioning (6 kVA and night time-of-use) and across all labels. It can be seen that samples generated by Latent Diffusion are close, in average, to the real test dataset, whereas samples from TimeGAN deviate slightly. Mean plots of other conditioning categories, not shown here, yield the same conclusions.

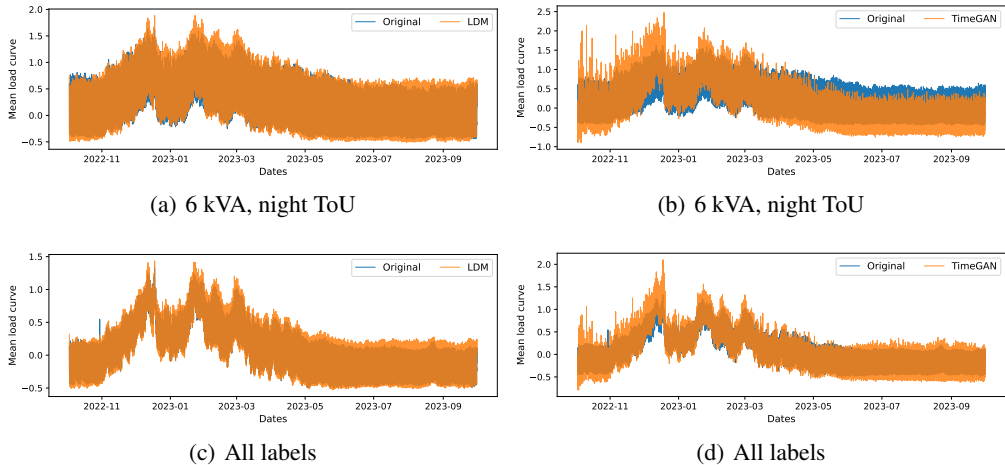


Figure 3: Average one-year load curves of (left) Latent Diffusion and (right) TimeGAN.

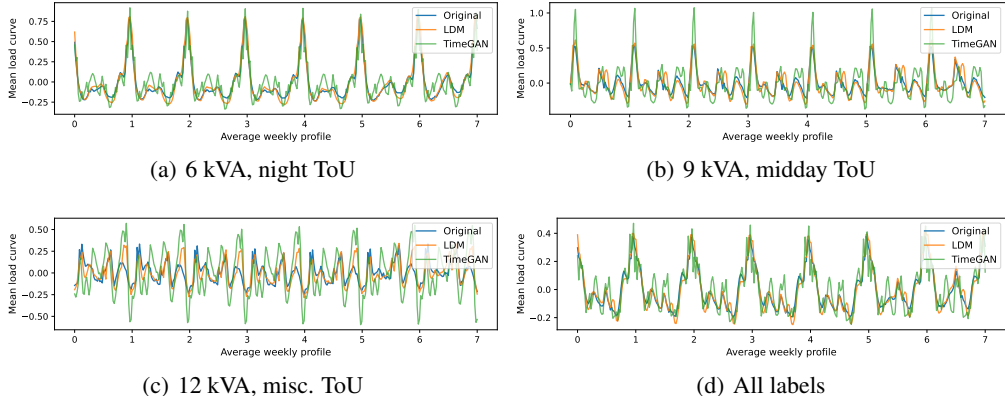


Figure 4: Average weekly profiles of Latent Diffusion (orange) and TimeGAN (green) against test data (blue).

Then, we compute the average weekly consumption of every sample, resulting in a 336-dimensional vector, and plot these mean profiles per label conditioning. A representative set of such profiles is shown in Figure 4. Different time-of-use rates induce strong weekly patterns, which are accurately reproduced by latent diffusion for the night and midday rates. This illustrates the flexibility of the conditioning mechanism, since label information is only included in the diffusion model and not in learning the autoencoder. On the contrary, TimeGAN’s profiles are noisier, although they capture the trend. We note that the *misc.* ToU patterns are harder to reproduce, possibly because they gather more diverse rate plans and because they represent about 10% of the training and test data.

As for temporal consistency, both Latent Diffusion and TimeGAN accurately capture the daily periodicity in electricity consumption (peak at lag 48), which is shown by the plots of the autocorrelation functions in Figure 5. However, TimeGAN adds spurious correlations at other lags, whereas Latent Diffusion shows high fidelity, both in mean and variance.

D.1.2 THERMO-SENSITIVITY

Gradients Figure 6 shows some thermo-sensitivity gradients – see Section A.3 – for the two models, compared to the hold-out test distribution. Load curves generated by Latent Diffusion have realistic gradients, matching real data. On the contrary, there is a slight mismatch on certain classes, e.g. for the category 6 kVA and night ToU, for TimeGAN.

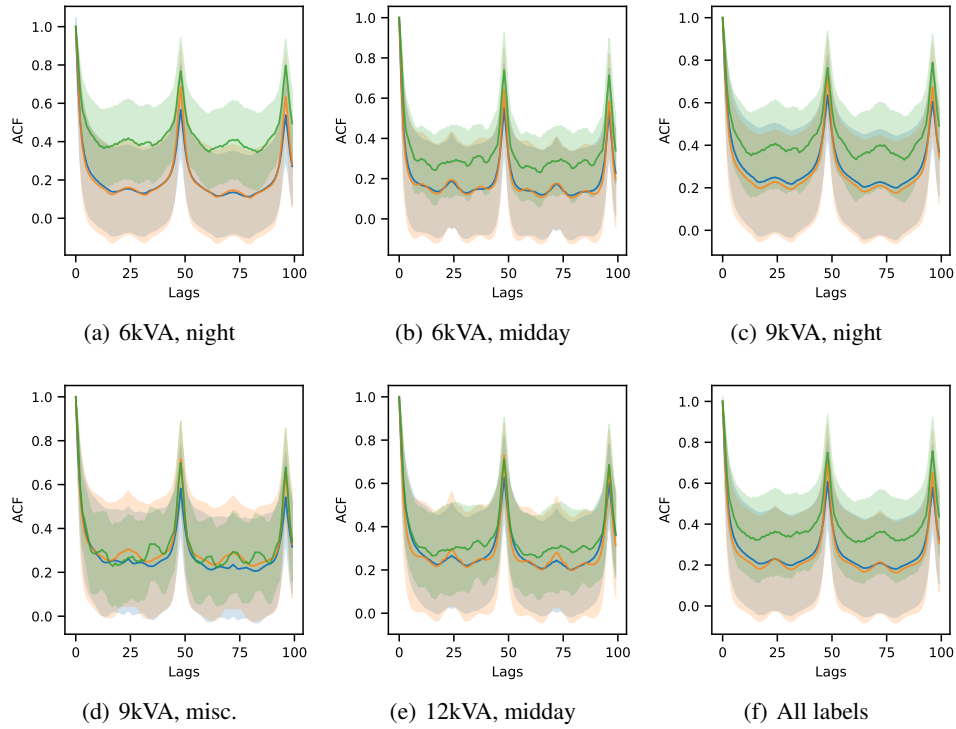


Figure 5: Average autocorrelation functions of Latent Diffusion (orange) and TimeGAN (green) against test data (blue). Shaded areas denote \pm standard deviation across samples.

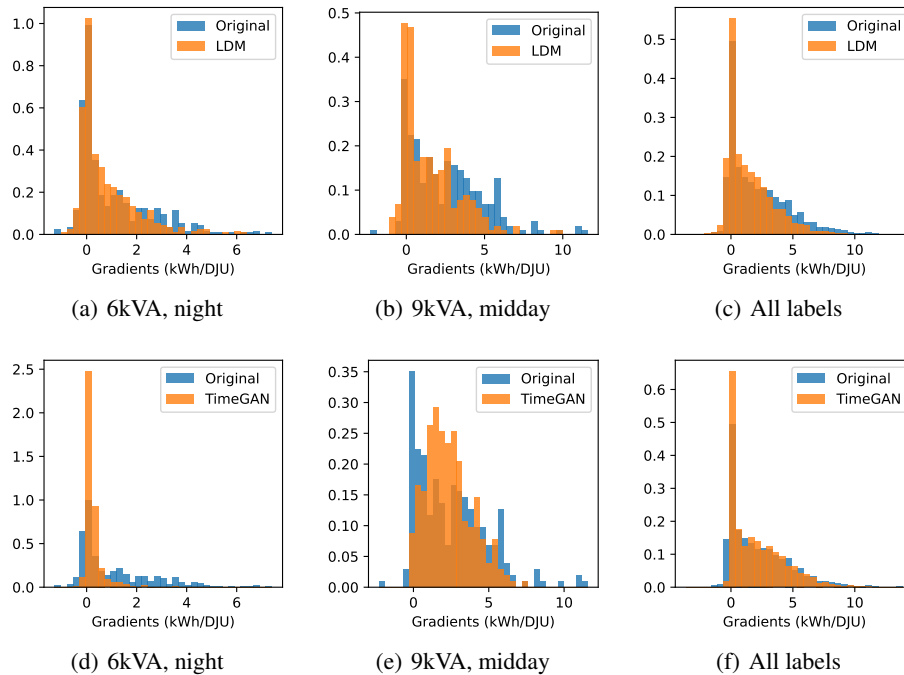


Figure 6: Thermo-sensitivity gradients computed from load curves aggregated at a daily frequency. Top: Latent Diffusion, bottom: TimeGAN.

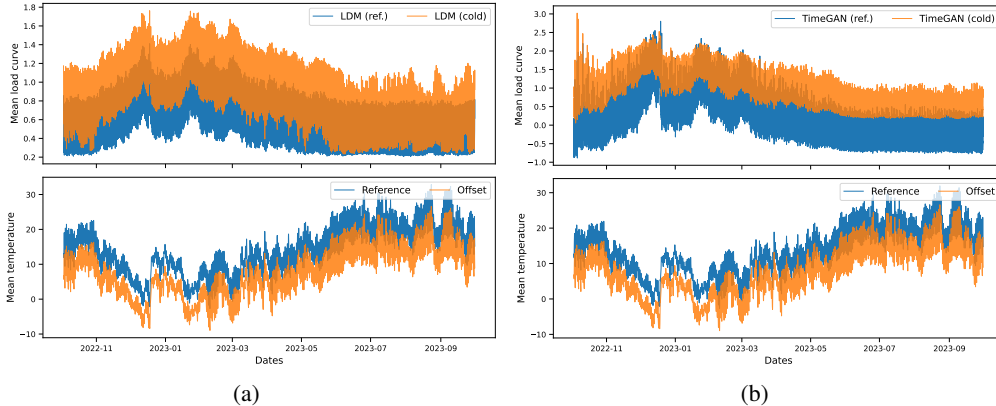


Figure 7: Average one-year synthetic load curves in the class 6 kVA, night ToU, with (orange) or without (blue) temperature offset for (a) Latent Diffusion and (b) TimeGAN.

Temperature offset Qualitatively, we assess how the synthetic load curves are distorted, in average, when generating under the same conditioning except that the temperature undergoes a strong -6.25°C offset throughout the year. We expect the consumption (i) to increase during winter days, since it is mostly driven by heating; (ii) to increase in summer, due to the strong offset, but less than in winter. Latent Diffusion samples in Figure 7 approximately follow this trend, particularly for winter days. There is a significant increase in summer as well, although the minimum values with or without offset are roughly the same. For TimeGAN, the strong temperature offset translates into a strong increase of the minimum electricity consumption all year long. However, the peak consumption is the same with or without offset during winter, which seems counter-intuitive if not unrealistic.

D.1.3 INDIVIDUAL SAMPLES

Finally, a random set of load curves generated by latent diffusion are shown in Figures 8, 9 and 10. Although it is challenging for the untrained eye to determine whether a given sample is realistic or not, it can be noted that the samples exhibit clear patterns, s.a. daily peaks, constant periods at low consumption, intermediate plateaus. The samples also show some diversity. To the best of our knowledge, these samples make thus viable load curves.

D.2 UTILITY

The full results for evaluating the forecasting utility are shown in Table 4. Load curves having known weekly seasonalities, we implement a simple baseline for the sake of comprehensiveness. To predict the next H values, the baseline repeats the last H values from one week ago, i.e. $\hat{\mathbf{x}}_{t+1:t+H} = \mathbf{x}_{t-336+1:t-336+H}$, for $H \leq 336$. Predicting individual load curves at sub-hourly frequency is notoriously challenging, particularly with no exogenous variables, due to the stochastic nature of user behaviors. This can be seen from the relatively modest improvement in mean absolute error (MAE) of PatchTST trained on Latent Diffusion over the baseline. Nevertheless, on both MSE and MAE and across all horizons, PatchTST trained on Latent Diffusion samples (i) is quasi-equivalent to its counterpart on real samples and (ii) outperforms the model trained on TimeGAN.

D.3 PRIVACY

MIA We plot the ROC curves in Figure 11 and report the black-box and white-box scores, namely the true positive ratio at 0.1% false positive ratio for membership inference attacks in Table 5, to summarize them (Carlini et al., 2022). Latent diffusion and TimeGAN obtain similar metrics: all scores are close to 0, typically around 0.002 (0.2%) for black-box and 0.1% for white-box, indicating a limited risk of re-identification. For the black-box attack, we also investigate whether the size of the available synthetic set, from which the smallest distance is computed, affects the scores. Except

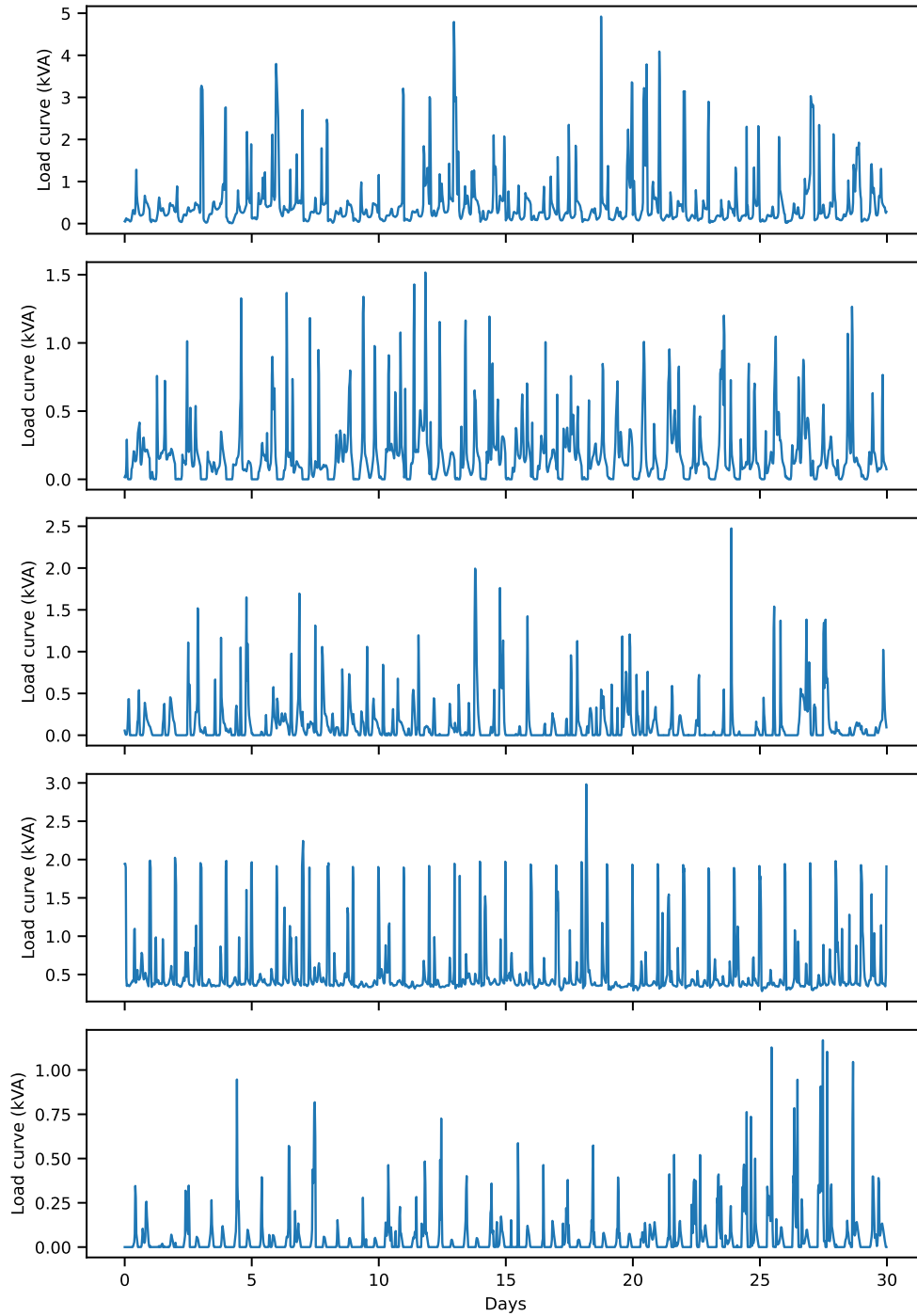


Figure 8: Individual samples (contracted power 6 kVA).

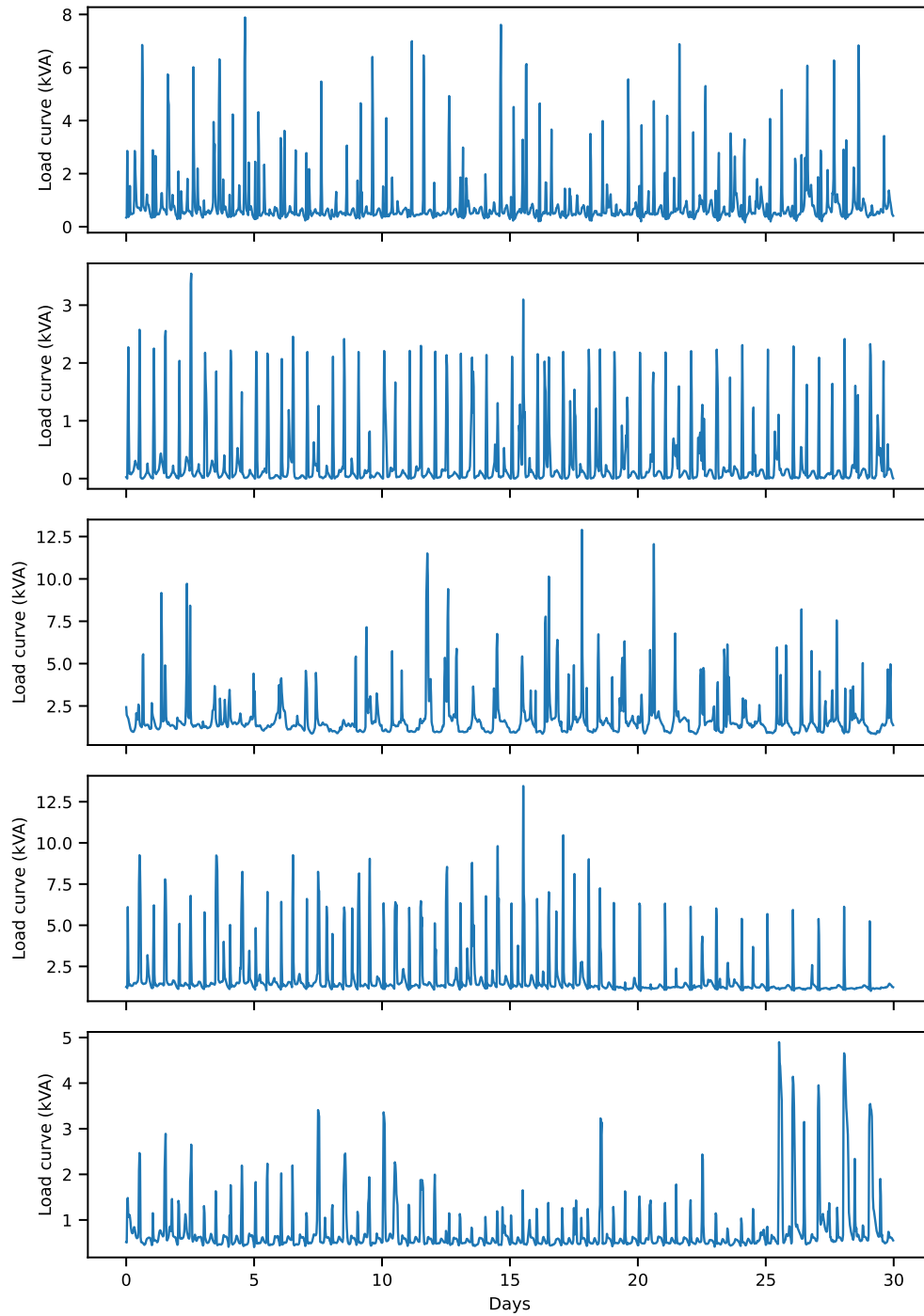


Figure 9: Individual samples (contracted power 9 kVA).

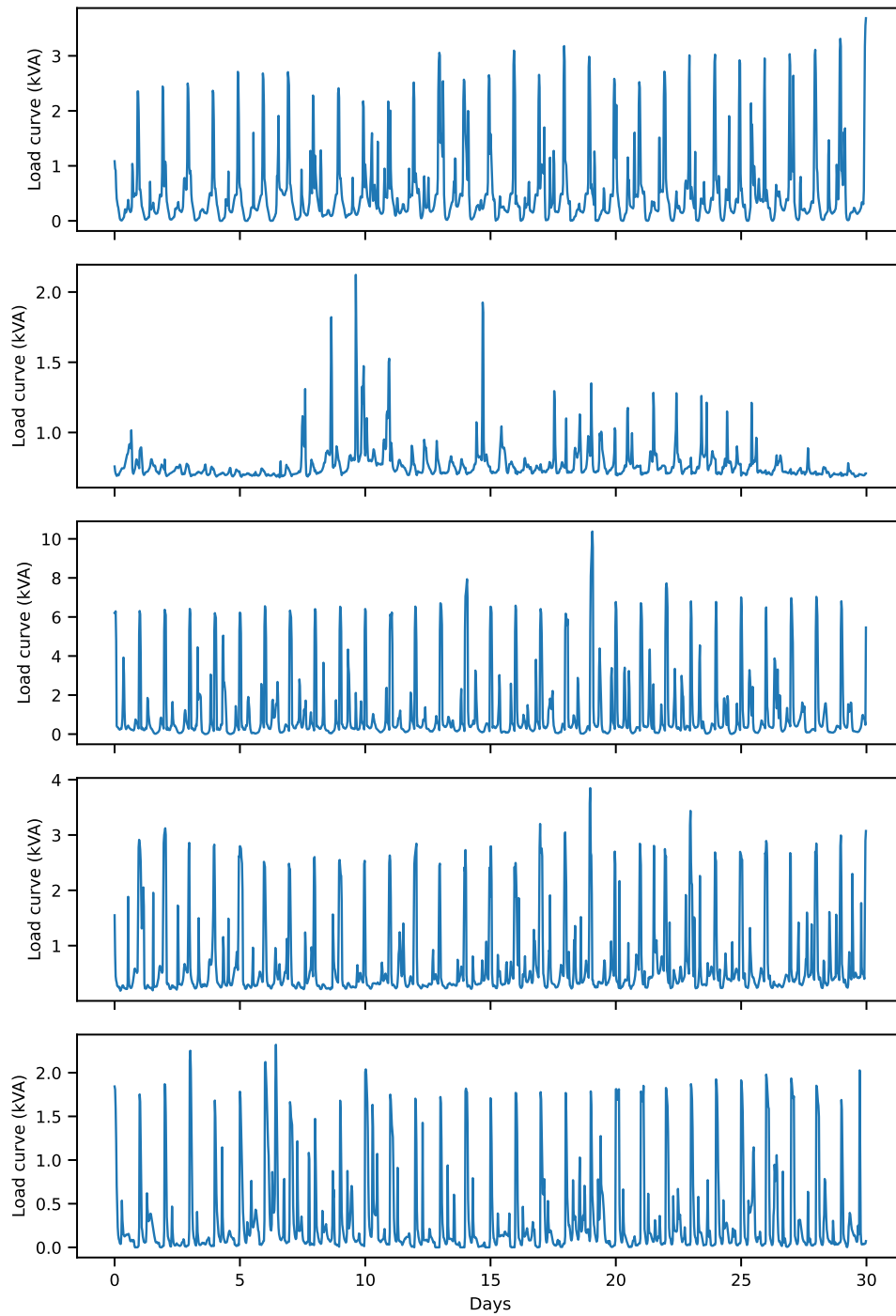


Figure 10: Individual samples (contracted power 12 kVA).

Table 4: Forecasting from a lookback of 720. Closest to TRTR (Train on Real, Test on Real) is emphasized in bold.

Horizon	Loss	TRTR	LDM	TimeGAN	Repeat
48	MSE	0.204	0.204	0.225	0.334
	MAE	0.240	0.237	0.259	0.264
96	MSE	0.188	0.188	0.207	0.308
	MAE	0.230	0.228	0.250	0.252
192	MSE	0.177	0.177	0.193	0.290
	MAE	0.228	0.227	0.245	0.243
336	MSE	0.192	0.192	0.211	0.294
	MAE	0.239	0.239	0.257	0.246

Table 5: True Positive Ratio ($\times 10^3$) at 0.1% False Positive Ratio for the membership inference attacks. The scores of the black-box attacks are averaged over 5 random subsamples (with \pm standard deviation) of size n_{gen} of synthetic samples.

n_{gen}	Black-box					White-box
	10	500	1k	1.5k	2k	-
LDM	0.43(± 0.36)	2.66(± 0.50)	2.64(± 0.37)	2.47(± 0.11)	2.26(± 0.13)	1.216
TimeGAN	1.19(± 0.74)	2.59(± 0.36)	2.52(± 0.11)	2.31(± 0.10)	2.21(± 0.03)	0.637

for very small sizes, Table 5 shows little variations in the scores, with a slight downward trend. Overall, the attacks on both LDM and TimeGAN are close to random and fail to reliably identify a subset of users.

Three-sample MMD Test The null hypothesis \mathcal{H}_0 is that the synthetic set of load curves is closer to the real test set than to the real train set. We compute the p-values of the three-sample MMD-test for all 9 categories of time-of-use \times contracted power, as well as for the entire synthetic dataset, in Table 6. The p-values are reasonably large for all categories and both Latent Diffusion and TimeGAN, suggesting not to reject \mathcal{H}_0 . Hence, this test does not give evidence that the models overfit the training data.

Nearest Neighbor Distance Ratio For both Latent Diffusion and TimeGAN, the distribution of NNDR computed with either real train or test data is the same, as evidenced in Figure 12. They

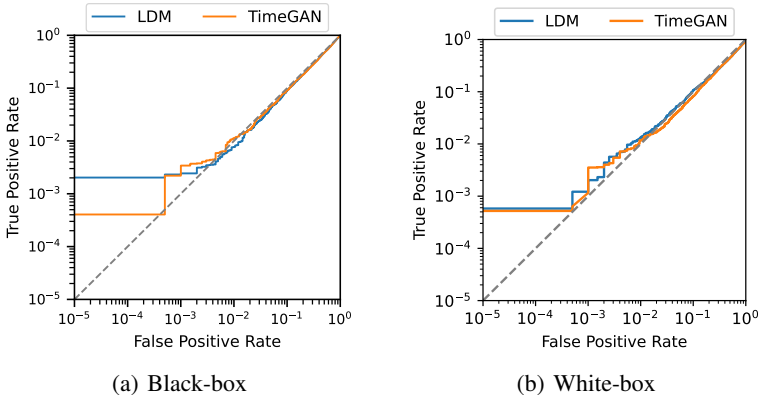


Figure 11: ROC curves in log scale of the membership inference attacks in the (a) black-box and (b) white-box scenarios.

Table 6: p-values for the MMD three-sample test (Bounliphone et al., 2016).

ToU	Midday			Night			Misc			all
Power (kVA)	6	9	12	6	9	12	6	9	12	all
LDM	0.297	0.930	0.852	0.019	0.218	0.795	0.239	0.325	0.642	0.845
TimeGAN	0.304	0.953	0.996	0.952	0.813	0.840	0.043	0.341	0.366	0.948

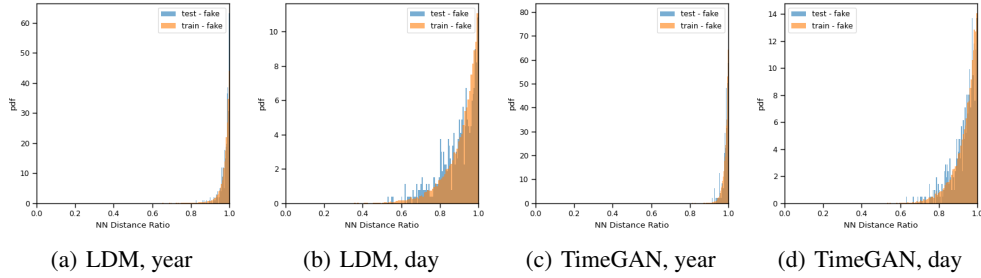


Figure 12: Nearest neighbor distance ratio for Latent Diffusion and TimeGAN, computed either on one-year data or average daily profiles.

are skewed towards 1, suggesting that the synthetic samples are not located towards outliers in the training set, which would constitute a breach of privacy.