
Operator Deep Smoothing for Implied Volatility

Lukas Gonon
Imperial College London
l.gonon@imperial.ac.uk

Antoine Jacquier
Imperial College London
a.jacquier@imperial.ac.uk

Ruben Wiedemann*
Imperial College London
r.wiedemann22@imperial.ac.uk

Abstract

We devise a novel method for implied volatility smoothing based on neural operators. The goal of implied volatility smoothing is to construct a smooth surface that links the collection of prices observed at a specific instant on a given option market. Such price data arises highly dynamically in ever-changing spatial configurations, which poses a major limitation to foundational machine learning approaches using classical neural networks. While large models in language and image processing deliver breakthrough results on vast corpora of raw data, in financial engineering the generalization from big historical datasets has been hindered by the need for considerable data pre-processing. In particular, implied volatility smoothing has remained an instance-by-instance, hands-on process both for neural network-based and traditional parametric strategies. Our general *operator deep smoothing* approach, instead, directly maps observed data to smoothed surfaces. We adapt the graph neural operator architecture to do so with high accuracy on ten years of raw intraday S&P 500 options data, using a single set of weights. The trained operator adheres to critical no-arbitrage constraints and is robust with respect to subsampling of inputs (occurring in practice in the context of outlier removal). We provide extensive historical benchmarks and showcase the generalization capability of our approach in a comparison with SVI, an industry standard parametrization for implied volatility. The operator deep smoothing approach thus opens up the use of neural networks on large historical datasets in financial engineering.

1 Introduction

Options trading experienced phenomenal growth in recent years. In its 2023 trading volume report [1], the CBOE announced the fourth consecutive year of record-breaking volumes on its options exchanges, citing a record-breaking number of transactions for European options on the S&P 500 index. European options are financial derivative contracts that give their holder the right, but not the obligation, to either buy or sell an underlying asset at a predetermined price (the *strike*) at a predetermined time (the *expiry*). An option specifying the right to buy (respectively to sell) is called a *Call* (respectively *Put*) option. Options are traded on a wide range of underlyings, including stocks, indices, currencies and commodities, and can be used to hedge against or speculate on the price movements of the underlying asset.

A key concept in options trading is the so-called *implied volatility*, which transforms the nominal price of an option into a conceptually and numerically convenient metric. The *implied volatility surface* is the collection of implied volatilities as observed at a specific point in time, visualized

*Corresponding author.

in three-dimensional space as a surface over the (strike, expiry)-domain. It provides an intuitive representation of the state of the options market and is crucial for hedging and risk management. The extraction of a smooth surface from quoted option prices is called *implied volatility smoothing* and remains one of the key challenges in options trading.

Conventionally, implied volatility smoothing relies on parametric surfaces whose parameters are optimized based on the distance to observed prices while adhering to *absence-of-arbitrage* conditions, which ensure the consistency of prices extrapolated from the smoothed surface. The development of such ad-hoc models for implied volatility traces back to *SVI* [10], which models implied volatility slice-wise for each maturity and successfully captures its key features on Equity indices. A continuous interpolation scheme for *SVI* slices was provided in [11], yielding a full surface. Nowadays, sophisticated market makers employ custom parametrizations, which can be considered proprietary trading secrets and reduce *SVI* to a benchmark role.

Regardless of the particular parametric surface model used, these approaches boil down to the continued execution of a numerical optimization routine: A smoothed surface expires as soon as quotes are updated (whenever markets move), necessitating the re-calibration of parameters. Success and duration of this routine is sensitive to initial conditions, search heuristics, and termination criteria, which exposes practitioners to considerable uncertainties during trading hours (or *online*). In response, we introduce a novel *operator deep smoothing* approach, replacing the instance-by-instance optimization with a single evaluation of a neural network. This greatly simplifies online calibration, at the upfront cost of training the network *offline* from historical data (in the spirit of [16, 17, 25]). Our unique use of *neural operators* [20] is fundamentally directed by the observation that the raw inputs for volatility smoothing (the collections of observed volatilities) vary in size and spatial arrangement: Options expire, new maturities and strikes become available, and the coordinates of existing options evolve continuously in the domain of the implied volatility surface (Figure 1b). This setting excludes classical neural networks – which required fixed-size inputs – from direct application. Neural operators, instead, conceptualize observed data as point-wise discretizations of latent functions in implicit infinite-dimensional function spaces and are well suited for the task. As the aim of volatility smoothing is the construction of a smooth surface that links the collection of market-observed option quotes, their use is highly motivated.

Contributions We introduce *operator deep smoothing*, a general approach for discretization-invariant data interpolation based on neural operators, and apply it to implied volatility smoothing. Our technique transcends traditional parametric smoothing and directly maps quoted volatilities to smoothed surfaces. Standard methods are limited to certain option markets (e.g. FX markets where options spread out on fixed rectilinear grids [5]) or require substantial data pre-processing. Instead, our technique adapts the graph neural operator (GNO) architecture [3] to consistently smooth input data of variable number and spatial arrangement. While neural operators have been successfully used in Physics to numerically solve partial differential equations, our application is the first in Finance and highlights the values of their discretization-invariance properties, so far rather under-explored. We test our method on ten years of intraday S&P 500 options data, accounting for more than 60 million volatility datapoints. Our method improves on the errors reported in [2], where much shorter periods of data were considered. In contrast, here we employ nine years of training data, and obtain double performance compared to the industry benchmark *SVI* on the final year of data. We summarize the impact of our contributions as follows:

- **Neural Operators in Financial Engineering** – Our GNO approach consistently maps quoted volatilities to smooth surfaces and is well suited to handle the highly dynamic nature of options data, removing the need for instance-by-instance processing of incoming data. Surfaces (or higher-dimensional equivalents) similar to the implied volatility surface are ubiquitous in quantitative finance, and we expect our technique to be transferable and to streamline and robustify algorithms’ data pipelines.
- **Neural Operators for Discretization-Invariant Interpolation** – Our operator deep smoothing approach leverages (a tweaked version of) GNOs for their abilities to interpolate irregular mesh geometries. This constitutes the first application of neural operators for interpolation/extrapolation tasks and paves the way to future research on the versatility of the discretization-invariance of neural operators in industrial applications characterized by dynamic and spatially irregular data.

Broader Impact Financial markets are centerpieces of modern economies. Operator deep smoothing, an algorithm to achieve robust and realistic generalization from market data, can help financial

intermediaries monitor and control financial risks better, thereby contributing to improved market efficiency. The beneficiaries include entities like pension funds, who use options to safeguard payment plans, and mutual funds, seeking to manage portfolio risks efficiently. Ultimately, operator deep smoothing may be useful for all market participants who own options in their portfolios. This includes the general public, whose trading participation has been increasing substantially [8]. Operator deep smoothing, providing "cheap" implied volatility smoothing with a single evaluation of a neural network, can thus help improve the quality of broadly accessible investment tools.

Literature Review The aforementioned SVI was developed for internal use at Merrill Lynch in 1999 and later advocated in [10]. Its extension to surface-based SSVI in [11] has been eagerly adopted by practitioners, which have since contributed to its robust calibration and generalizations [7, 15, 12]. It was augmented in [2] by a multiplicative neural network corrector, based on guided network training by means of no-arbitrage soft constraints from [34]. The absence-of-arbitrage conditions – providing safeguards for option pricing – for implied volatility surfaces were formulated in [31], and we provide an equivalent formulation, based on [9, 26], for practical purposes. In [6] static arbitrage constraints were used to perform option calibration (with an additional regularization technique), which can be considered to be instance-by-instance smoothing of nominal price data. In [5] a classical VAE (variational autoencoder) was applied to implied volatility smoothing on FX markets, where strikes of quoted options are tied to a fixed grid of *deltas*.² This specificity of FX markets allows the use of a conventional feedforward neural network based decoder. Recent option calibration approaches based on neural networks have been proposed in [4, 16, 17, 33].

A comprehensive account on neural operators is given in [20], unifying previous research on different neural operator architectures and techniques [3, 22]. Subsequent developments investigating the expressivity of these architectures as well as their generalizations include [13, 18, 21, 23, 24, 32].

Outline We review financial concepts and the challenges of implied volatility smoothing in Section 2. In Section 3, we provide a review of neural operators (Section 3.1) and introduce (Section 3.2) our operator deep smoothing approach for general interpolation tasks. In Section 4, we perform experiments for implied volatility smoothing of S&P 500 options data. Finally, Section 5 gathers limitations as well as outlooks regarding the use of neural operators for interpolation purposes.

Code We make code for the paper available at the location <https://github.com/rwic1/operator-deep-smoothing-for-implied-volatility>. In particular, the code repository contains a general *PyTorch* [30] implementation of the graph neural operator architecture for operator deep smoothing.

2 Background: Implied Volatility

We consider a market of European options written on an underlying asset, which we observe at a given instant T_0 and denote the time- T forward price of the underlying asset by $F_{T_0, T}$.

European Call Options The option market consists of a finite collection of *European Call options*,³ each identified by its expiry $T \in (T_0, \infty)$ and its strike $K \in (0, \infty)$, and we write $C(T, K)$ for its (undiscounted) price. In practice, these are traded for fixed expiries T_1, \dots, T_m ; for each T_i , only a finite range of strikes $K_1^i, \dots, K_{n_i}^i$ is available, typically widening as the expiry increases (Figure 1).

Black-Scholes The *Black-Scholes model* is the simplest diffusive asset model and captures the volatility of the underlying asset with a single parameter $v \in (0, \infty)$. Its popularity stems from the closed-form expression it admits for the price of a European Call option with *time-to-expiry* $\tau = T - T_0$ and *log-moneyness* $k = \log(K/F_{T_0, T})$:

$$\text{BS}(\tau, k, v) = \Phi(d_1(\tau, k, v)) - e^k \Phi(d_2(\tau, k, v))$$

²Delta is the derivative of the price with respect to the underlying asset and is standard in FX strike quoting.

³In practice, market participants trade both Call and Put options, which are mathematically equivalent through the well-known *Put-Call parity*. The latter thus allows to speak in terms of Call options only.

(in units of the time- T forward of the underlying). Here, Φ denotes the cumulative distribution function of the standard Normal distribution, while

$$d_1(\tau, k, v) = \frac{-k}{v\sqrt{\tau}} + \frac{1}{2}v\sqrt{\tau}, \quad d_2(\tau, k, v) = \frac{-k}{v\sqrt{\tau}} - \frac{1}{2}v\sqrt{\tau}. \quad (1)$$

This model also provides a closed-form expression for the sensitivity of the option price with respect to its volatility parameter. Termed *Vega* and required for our experiments in Section 4.1, it reads

$$\mathcal{V}(\tau, x, v) = \partial_v \text{BS}(\tau, k, v) = \varphi(d_1(\tau, k, v))\sqrt{\tau}. \quad (2)$$

Implied Volatility While not able (any longer) to fit market data, the mathematical tractability of the Black-Scholes model gave rise to the important concept of *implied volatility*: Given a Call option with price $C(T, K)$, its implied volatility $v(\tau, k)$ solves $C(T, K) = F_{T_0, T} \text{BS}(\tau, k, v(\tau, k))$. By using time-to-expiry/log-moneyness coordinates, the implied volatility provides a universal way to consistently compare the relative expensiveness of options of across different strikes, maturities, underlyings, and interest rate environments. Its characteristic shape helps traders make intuitive sense of the instantaneous state of option markets relative to a (flat) Black-Scholes model baseline.

Implied Volatility Smoothing This refers to fitting a smooth surface $\hat{v}: (0, \infty) \times \mathbb{R} \rightarrow (0, \infty)$ to a collection $\mathbf{v} = \{v(\tau_l, k_l)\}_{l=1}^p$ of observed implied volatilities. Naive strategies such as cubic interpolation, are ill-fated as they yield *arbitrage*, namely cost-less trading strategies generating a guaranteed profit. In option markets, an arbitrage is called *static* when set up solely from fixed positions in options and a dynamically (but a finite number of times) readjusted position in the underlying. Beyond simple interpolations, practitioners have devised ad-hoc parametrizations for implied volatility, in particular the aforementioned SVI, which are not expected to perfectly match all reference prices. Instead, the model parameters are optimized with respect to an objective function that measures market price discrepancy and includes penalization terms ruling out static arbitrage. These penalization terms are commonly formulated on the basis of the following theorem, which summarizes the shape constraints of the implied volatility surface [11, 26, 31].

Theorem 2.1 (Volatility Validation). *Let $\hat{v}: (0, \infty) \times \mathbb{R} \rightarrow (0, \infty)$ be continuous and satisfying*

- (i) *Calendar arbitrage: For each $k \in \mathbb{R}$, $\hat{v}(\cdot, k)\sqrt{\cdot}$ is non-decreasing and vanishes at the origin.*
- (ii) *Strike arbitrage: For every $\tau > 0$, the slice $\hat{v}_\tau = \hat{v}(\tau, \cdot)$ is of class C^2 with*

$$\text{But}(\tau, \cdot, \hat{v}_\tau, \partial_k \hat{v}_\tau, \partial_k^2 \hat{v}_\tau) \geq 0 \quad (3)$$

and $\limsup_{k \uparrow \infty} \frac{\hat{v}_\tau^2(k)}{k} < \frac{2}{\tau}$, where

$$\text{But}(\tau, k, v_0, v_1, v_2) = (1 + d_1(\tau, k, v_0)v_1\sqrt{\tau})(1 + d_2(\tau, k, v_0)v_1\sqrt{\tau}) + v_0v_2\tau.$$

Then, $(T, K) \mapsto \text{BS}(\tau, k, \hat{v}(\tau, k))$ defines a Call price surface that is free from static arbitrage.

Condition 2.1(i) is equivalent to prices increasing in maturity (uncertain increases as time passes), while Condition 2.1(ii) arises when computing the *implied probability density* f_τ of the underlying:

$$f_\tau(\cdot) = \frac{\varphi(-d_2(\tau, \cdot, v))}{v} \text{But}(\tau, \cdot, v, \partial_k v, \partial_k^2 v). \quad (4)$$

Since a density needs to be non-negative, Identity (4) explains why Condition 2.1(ii) above is required.

3 Neural Operators for Discretization-Invariant Smoothing

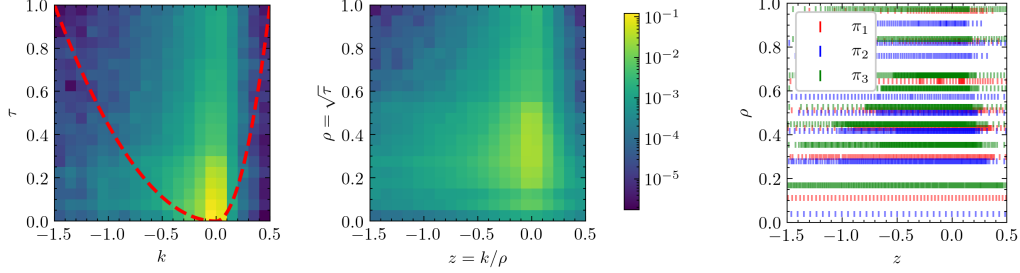
3.1 Background: Neural Operators

We provide full details about notations, terms, and additional context in Appendix A.

Philosophy The development of *neural operators* is based on the philosophy that observed data $\mathbf{a} = \{a_l\}_{l=1}^p$ arises as the evaluation of a latent function $a: D \rightarrow \mathbb{R}^{c_n}$, defined on some domain $D \subseteq \mathbb{R}^d$, at a discretization $\pi = \{x_l\}_{l=1}^p$ of D . That is, $\mathbf{a} = a|_\pi$, or

$$a_l = a(x_l), \quad l = 1, \dots, p. \quad (5)$$

An input-output relationship of data $\mathbf{a} \mapsto \mathbf{u}$ is then "really" described by an operator $F: \mathcal{A} \rightarrow \mathcal{U}$ between function spaces \mathcal{A} and \mathcal{U} . Neural operators are abstract neural network architectures $F^\theta: \mathcal{A} \rightarrow \mathcal{U}$, with implementations that integrate (5) consistently across the variable discretization π .



(a) Relative trading volume per quoted interval, averaged over S&P 500 dataset 2012-2021. Left: Rectangular domain w.r.t. time-to-expiry/log-moneyness (99.9% of trading volume with maturities < 1 year). Right: Rectangular domain w.r.t. transformed coordinates (96.6% of trading volume while more evenly populated); boundary delineated in Left. (b) Scatter plots of example sets of option quotes $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3$. Observed on 15.10.2012, 18.05.2017, and 04.01.2021, respectively, at 10:50:00, each.

Figure 1: Spatial arrangement of option quotes in time-to-expiry/log-moneyness domain.

Technicalities We review the core concepts of neural operators from [20]. Let D be a bounded domain in \mathbb{R}^d and \mathcal{A} and \mathcal{U} be Banach spaces of functions mapping from D to $\mathbb{R}^{c_{\text{in}}}$ and $\mathbb{R}^{c_{\text{out}}}$, respectively. Neural operators are finitely parametrized mappings $F^\theta: \mathcal{A} \rightarrow \mathcal{U}$ with *universality* for continuous target operators and with *discretization-invariant* implementations \tilde{F}^θ . In the space $C(\mathcal{A}, \mathcal{U})$ topologized by uniform convergence on compacts, the architecture F^θ is called *universal* if $\{F^\theta\}_{\theta \in \Theta}$ is dense in $C(\mathcal{A}, \mathcal{U})$, with Θ the parameter set. An implementation of F^θ is an algorithm \tilde{F}^θ which accepts observed data $\mathbf{a} = a|_\pi$ and outputs a function $u \in \mathcal{U}$ and is such that $\tilde{F}^\theta_\pi(\cdot) = \tilde{F}^\theta(\cdot|_\pi) \in C(\mathcal{A}, \mathcal{U})$. Now, \tilde{F}^θ is called *discretization-invariant* if $\lim_{n \uparrow \infty} \tilde{F}^\theta_{\pi^{(n)}} = F^\theta$ in $C(\mathcal{A}, \mathcal{U})$, given a *discrete refinement*⁴ of D .

Let K be a set of input functions, compact in \mathcal{A} , and let $\varepsilon > 0$. In combination, universality and discretization-invariance allow to posit the existence of parameters θ , such that for all $a \in K$,

$$\|\tilde{F}^\theta(a|_\pi) - F(a)\|_{\mathcal{U}} \leq \varepsilon, \quad (6)$$

irrespective of the particular discretization π given that the data $\mathbf{a} = a|_\pi$ is scattered sufficiently densely across D . The training of neural operators is analogous to the classical finite-dimensional setting. It happens in the context of an implicit *training distribution* μ on the input space \mathcal{A} and aims at minimizing the *generalization error*

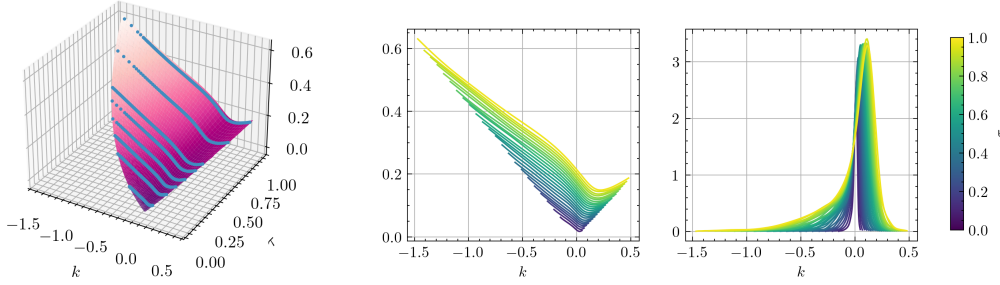
$$R_\mu: \theta \mapsto \mathbb{E}_{a \sim \mu} \|\tilde{F}^\theta(a|_\pi) - F(a)\|_{\mathcal{U}}, \quad (7)$$

through the use of gradient descent methods applied to empirical estimates of (7). These estimates are constructed from a *training dataset* $\mathcal{D} = \{(\mathbf{a}^{(i)}, \mathbf{u}^{(i)})\}_{i=1}^n$ of *features* $\mathbf{a}^{(i)} = a^{(i)}|_{\pi^{(i)}}$ and *labels* $\mathbf{u}^{(i)} = F(a^{(i)})|_{\pi^{(i)}}$ on the basis of (*mini*) *batching* heuristics and are frequently transformed or augmented by additional terms through the use of apposite *loss functions*.

3.2 Operator Deep Smoothing

Let $\mathbf{v} = \{v(x_l)\}_{l=1}^p$ be the collection of observed data, for example implied volatilities as in Section 2. This notation silently adopts the neural operator philosophy, connecting the data point $v(x_l)$ with coordinates x_l , hinting at a latent function $v: D \rightarrow \mathbb{R}$ giving rise to the observed values. The *smoothing* or *interpolation* task consists of constructing an appropriately regular candidate \hat{v} for v from the data $\mathbf{v} = v|_\pi$, with $\pi = \{x_1, \dots, x_p\}$. The operator deep smoothing approach uses a neural operator \tilde{F}^θ , trained using historical data, to generate $\hat{v} := \tilde{F}^\theta(\mathbf{v})$. It fundamentally leverages the discretization-invariance to produce consistent results even if the "sensors" x_l of the latent function v continuously change in availability and/or location in the domain D . This is the situation for volatility smoothing, where the sensors $x_l = (\tau_l, k_l)$ are the (time-to-expiry, log-moneyness) coordinates of the quoted options, which – as noted in Section 1 and illustrated in Figure 1b – move with the market,

⁴A *discrete refinement* of D is a nested sequence $(\pi^{(n)})_{n \in \mathbb{N}}$ of discretizations of D for which for every $\varepsilon > 0$, there exists $N \in \mathbb{N}$ such that $\{\mathbb{B}_{\mathbb{R}^d}(x, \varepsilon) : x \in \pi^{(N)}\}$ covers D .



(a) Smoothed surface with input volatility quotes (blue). $\delta_{\text{abs}} = 0.004$, $\delta_{\text{spr}}(\hat{v}, \mathbf{v}) = 1.061$. (b) Slices of implied volatility (left) and implied density (right), scaled by time-to-expiry. Left: Absence of crossings indicates absence of calendar arbitrage. Right: Positivity indicates absence of butterfly arbitrage. In fact $\mathcal{L}_{\text{cal}}(\theta, \mathbf{v}) = 6.3 \times 10^{-6}$, $\mathcal{L}_{\text{but}}(\theta, \mathbf{v}) = 0$.

Figure 2: Operator deep smoothing of quotes \mathbf{v} from 04.01.2021 at 10:50:00. Compare Figure 10.

thus motivating the operator deep smoothing approach. We now describe the methodology in more detail. For a first glance at the results we refer to Figure 2.

Methodology Translating the smoothing task to an operator learning problem, we find that the target operator is the (continuous) identity operator $F = F_{\text{id}}: \mathcal{A} \ni v \mapsto v \in \mathcal{U}$, with equal input and output spaces $\mathcal{A} = \mathcal{U} = C(\bar{D})$. The training dataset is the collection $\mathcal{D} = \{\mathbf{v}^{(i)}\}$ of historical data (labels and features coincide), and we train a suitable neural operator architecture \tilde{F}^θ such that

$$|\tilde{F}^\theta(\mathbf{v})(x) - v(x)| \leq \varepsilon, \quad \text{for all } x \in D, \quad (8)$$

for some given error tolerance ε for the interpolation task, while v is the latent function of which we observe $\mathbf{v} = v|_\pi$. Instead of minimizing empirical estimates of the ∞ -norm $\|\tilde{F}^\theta(\mathbf{v}) - v\|_{\mathcal{U}}$, however, we suggest a *fitting loss* based on the *root mean square relative error*,

$$\mathcal{L}_{\text{fit}}(\theta, \mathbf{v}) := \sqrt{\frac{1}{|\pi|} \sum_{x \in \pi} \left(\frac{|\tilde{F}^\theta(\mathbf{v})(x) - v(x)|}{1 \vee |v(x)|} \right)^2}, \quad (9)$$

for its smoothness and invariance to the scale of the data.⁵ Additional engineering techniques, such as sub-sampling of inputs during training, are explored in our practical investigation in Section 4.

Practical Constraints Depending on the application, the smoothing task may be subject to constraints. For volatility smoothing, the smoothed surface $\hat{v}^\theta = \tilde{F}^\theta(\mathbf{v})$ must be free of static arbitrage. This is effectively enforced by augmenting the loss function with additional penalization terms, moving away from a pure operator learning problem.⁶ This does not only promote the relevant properties in the neural operator output but can also help define it when faced with sparsity of data in the domain D (in this context see also [23]). From Theorem 2.1, these penalization terms naturally motivate our choices of \mathcal{L}_{but} and \mathcal{L}_{cal} . The strike arbitrage constraint 2.1(ii) is handled via

$$\mathcal{L}_{\text{but}}(\theta; \mathbf{v}) = \left\| \left(\text{But}(\cdot, \hat{v}^\theta, \partial_k \hat{v}^\theta, \partial_k^2 \hat{v}^\theta) - \varepsilon \right)^- \right\|_1, \quad (10)$$

where we ignore the asymptotic condition since our experiments are focused on the bounded domain D (Figure 1a). The inclusion of ε promotes strictly positive implied densities (we will use $\varepsilon = 10^{-3}$), while we choose the 1-norm to induce sparsity in the constraint violation. The calendar arbitrage constraint 2.1(i) can be tackled analogously with

$$\mathcal{L}_{\text{cal}}(\theta; \mathbf{v}) = \left\| \left(\partial_\tau [(\tau, k) \mapsto v^\theta(\tau, k)\sqrt{\tau}] - \varepsilon \right)^- \right\|_1, \quad (11)$$

where again we ignore the asymptotic condition since D is bounded away from zero time-to-expiry.

⁵Empirical estimates of ∞ -norms and L^2 -norms are equivalent loss functions on finite-dimensional spaces.

⁶It is not the goal of volatility smoothing to learn the identity operator throughout the entire input space.

Interpolating Neural Operator Architectures Various neural operator architectures exist, mostly arising from the *kernel integral transform* framework of [20]. Most prominently, these include Fourier neural operators (FNO), delivering state-of-the-art results on fixed grid data, as well as graph neural operators (GNO), able to handle arbitrary mesh geometries, both reviewed in Appendix A.1. While highly effective with documented universality, these neural operators are not directly applicable for interpolation as their layers include a pointwise-applied linear transformation, which limits the output to the set of the input data locations. Dropping this *local* linear transformation results in an architecture proved to retain universality [20] and – at least for its implementation as a GNO – fit for interpolation tasks. In fact, universality holds [21] for an architecture combining the local linear transformation with a simple averaging operation, suggesting the fundamental importance of the collaboration of local and non-local components; this was noted in [20] for whom retaining the local components can be "beneficial in practice". We confirm this in our experiments showing that a purely non-local architecture has substantially reduced performance.

We therefore propose a new architecture for operator deep smoothing leveraging GNOs' unique ability to handle irregular mesh geometries. We use a purely non-local first layer (dropping the pointwise linear transformation), and use it to produce hidden states at all required output locations, enabling subsequent layers to retain their local transformations. Since GNOs do not theoretically guarantee a smooth output, we augment the training with additional regularization terms such as $\mathcal{L}_{\text{reg}}(\theta; \mathbf{v}) = \|\Delta \hat{v}^\theta\|_2$, with Δ the Laplace operator, and provide a full description in Appendix B.

4 Experiments

We detail our practical investigation of the operator deep smoothing approach for implied volatility.

4.1 Model Training

Dataset and Splits We perform our numerical experiments using 20-minute intervals of *CBOE S&P 500 Index Option* data from 2012 to 2021. The dataset amounts to a collection of 49089 implied volatility surfaces and just above 60 million individual volatility data points (after domain truncation). We refer the reader to Appendix C.1 for full details on the preparation of the dataset. We allocate the first nine years of data (2012 to 2020) to training, keeping 750 randomly drawn surfaces for validation purposes, and use the final year of the dataset (2021) for testing. This yields a training dataset $\mathcal{D}_{\text{train}}$ containing $n_{\text{train}} = 43442$ surfaces, a validation dataset \mathcal{D}_{val} containing $n_{\text{val}} = 750$ surfaces and a test dataset $\mathcal{D}_{\text{test}}$ with $n_{\text{test}} = 4897$ surfaces.

Data Transformation Motivated by Figure 1a, we transform the time-to-expiry and log-moneyness coordinates via $\rho = \sqrt{\tau}$ and $z = k/\rho$. Intuitively, this transformation converts the "natural" scaling of implied volatility by time-to-expiry to a scaling of the input domain. From here on, we consider the domain in these coordinates, setting $D = (\rho_{\min}, \rho_{\max}) \times (z_{\min}, z_{\max}) = (0.01, 1) \times (-1.5, 0.5)$. In (τ, k) -coordinates, D becomes a cone-shaped region, that, on average, contains 96.6% of traded options (with time-to-expiry below one year) and, with respect to (ρ, z) -coordinates, is more evenly populated, improving the numerics.

Model Configuration The model hyperparameters (giving rise to 102529 trainable parameters in total), identified by manual experimentation, are detailed in Appendix C.2.

Loss Function We implement a *Vega*-weighted (see (2)) version of the fitting loss \mathcal{L}_{fit} (9). We compute \mathcal{L}_{but} directly as (10) (in the transformed coordinates) on a synthetic grid using finite differences. For \mathcal{L}_{cal} , we implement a multiplicative version that is invariant to the scale of implied volatility. We provide a precise description in Appendix C.3.

Model Training We train the GNO F^θ over 500 epochs on the training dataset using the AdamW optimizer with learning rate $\lambda = 10^{-4}$ and weight decay rate $\beta = 10^{-5}$, and use a pseudo batch size of 64 by accumulating gradients. We randomly sub-sample the inputs \mathbf{v} and randomize the grids on which we compute the arbitrage losses. The training is performed in around 200 hours using a NVIDIA Quadro RTX 6000 GPU. The validation loss is reported in Appendix C.3 in Table 1.

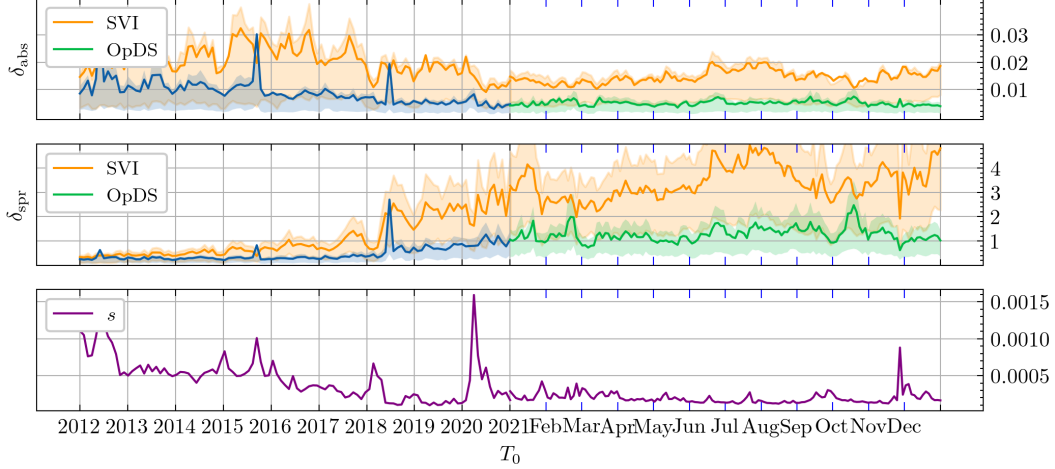


Figure 3: Benchmark metrics (surface-averages, between 25%- and 75%-quantiles) over training period (computed from \mathcal{D}_{val}) and testing period, as well as surface-averaged spread. Resampled, monthly for training and weekly for testing period.

4.2 Results

Benchmark Metrics Let $\mathbf{v} = \{v(x)\}_{x \in \pi}$ be the collection of observed implied volatilities and \hat{v} the smoothed surface as produced by a given method. We measure absolute relative error:

$$\delta_{\text{abs}}(\hat{v}(x), v(x)) = \frac{|\hat{v}(x) - v(x)|}{v(x)}.$$

The average of $\delta_{\text{abs}}(\hat{v}, \mathbf{v})$ is known as the *mean absolute percentage error* (or *MAPE*). As in [7], we moreover realize the importance of analyzing the smoothing algorithm in terms of nominal price error relative to the size of the bid-ask spread $s(x) = \text{BS}_{\pm}(x, v_{\text{ask}}(x)) - \text{BS}_{\pm}(x, v_{\text{bid}}(x))$.⁷ We define

$$\delta_{\text{spr}}(\hat{v}(x), v(x)) = \frac{2}{s(x)} |\text{BS}_{\pm}(x, \hat{v}(x)) - \text{BS}_{\pm}(x, v(x))|.$$

Since we use simple mid reference prices, $\delta_{\text{spr}}(\hat{v}(x), v(x)) \leq 1$ indicates that the prediction $\hat{v}(x)$ for the option x lies within the bid-ask spread.

Model Finetuning During production use, the GNO would be retrained regularly using the most recent available data. We emulate this procedure by evaluating the benchmark metrics on a monthly basis over the 2021 test dataset $\mathcal{D}_{\text{test}}$: Following the evaluation of the first month’s test data, the GNO is trained for 10 epochs on the test data just reviewed, with each mini-batch augmented by an equal amount of data from the training dataset $\mathcal{D}_{\text{train}}$. We repeat this, progressively incorporating an additional month of data, until the entire dataset $\mathcal{D}_{\text{test}}$ is assessed. This finetuning-evaluation procedure takes circa 1.8 GPU hours per month.

Analysis It is apparent from Figure 3 that the operator deep smoothing approach substantially improves on SVI’s smoothing capabilities, with respect to both δ_{abs} and δ_{spr} .⁸ The use of our approach in practice, with monthly finetuning, can be expected to smooth the volatility surface with a MAPE of around 0.5%, while SVI fluctuates between 1% and 2%. Our approach is competitive with [2], which perform instance-by-instance volatility smoothing using neural networks, and report a MAPE

⁷ $v_{\text{bid}}(x)$ and $v_{\text{ask}}(x)$ are the implied volatilities corresponding to Bid and Ask option prices while BS_{\pm} is the Black-Scholes formula for Call (resp. Put) options for positive (resp. negative) log-moneyness values:

$$\text{BS}_{\pm}(\tau, k, v) = \begin{cases} \Phi(d_1(\tau, k, v)) - e^k \Phi(d_2(\tau, k, v)), & k > 0 \\ e^k \Phi(-d_2(\tau, k, v)) - \Phi(-d_1(\tau, k, v)), & k \leq 0 \end{cases}$$

⁸We produce the SVI benchmark as described in Section E.

of around 1% for the period Jan-Apr 2018.⁹ The various figures in Appendix D show the qualitative improvements of our method versus the SVI benchmark.

An important metric to consider when making sense of Figure 3 is the historical tightening of the bid-ask spread, driven by increases in competition on the S&P 500 option market over recent years. This fact explains why δ_{spr} is very small on the chronologically earlier training dataset, both for operator deep smoothing and SVI, even as δ_{abs} is comparatively larger. The latter fact is explained similarly (and is visually recognizable by Figure 8): Wide spreads give rise to more variable prices, leading to greater need for correction by the smoothing algorithm, captured by δ_{abs} .

Complementary to Figure 3, Figure 4 resolves the error metrics as well as the terms controlling the absence of arbitrage spatially, averaged along the time axis. The MAPE tends to be larger on the Call wing, in accordance with noisier option prices (Call options experience less trading than Put options). Moreover, we discern that, on average, the smoothed surfaces are completely free of arbitrage.

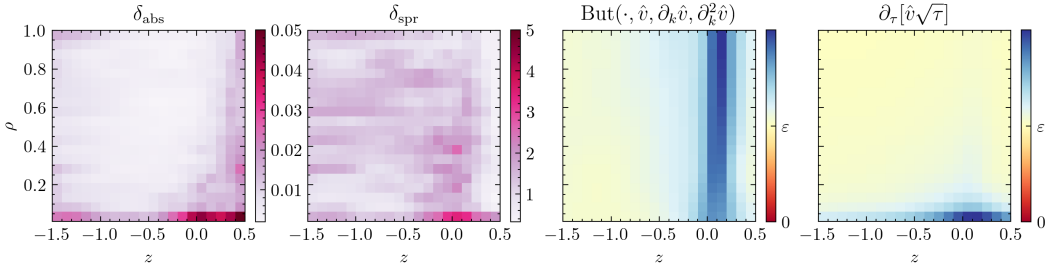


Figure 4: Average spatial distribution of benchmark metrics and arbitrage terms over $\mathcal{D}_{\text{test}}$.

5 Discussion

Summary We provide a novel method for implied volatility smoothing, resulting from an application of our general operator deep smoothing approach for discretization-invariant data interpolation. The approach leverages a graph neural operator to directly map given data – consistently across data dimensions and spatial arrangement – to smoothed surfaces, transcending classical parametric smoothing techniques. In the example of volatility smoothing, benefits include a massively simplified online calibration process.

Subsampling of Inputs Due to the discretization-invariance of the GNO architecture underlying the operator deep smoothing approach, the smoothing method is robust with respect to subsampling of inputs (we include this technique into operator training to improve generalization, compare Appendix C.3). In practice, subsampling of inputs occurs in the context of outlier removal. In the example of volatility smoothing, certain quotes may be determined spurious. Simply removing anomalous datapoints from the input is compatible with our method.

Compression Figure 3 makes the compression qualities of the operator deep smoothing approach apparent: We compute the entire historical timeseries using a single GNO instance, with around 100 thousand parameters. Evaluating the SVI benchmark, on the other hand, requires 61454 model instances (one per slice), or a total of 307270 parameters. A comparison with [2], which for each smoothed surface instantiates a deep neural network corrector of around 5085 parameters,¹⁰ is striking. Performing [2]’s suggested smoothing on the entire dataset 2012–2021 would require more than 200 million parameters. On the other hand, we expect our operator deep smoothing approach to deliver competitive results continuously over the entire training period and beyond (with regular finetuning), even when moving to higher-frequency data (our data is summarized on a 20-minute interval basis).

⁹However, [2] does not perform a similar restriction of the domain of the volatility surface.

¹⁰Computed as the sum of $120 = 3 \times 40$ parameters for the input layer, three times $1640 = 41 \times 40$ parameters for the hidden layers, 41 parameters for the output layers, plus 4 additional parameters of the SSVI prior and a scaling parameter.

Limitations and Perspectives Compared to ad-hoc volatility parametrizations like SVI, the operator deep smoothing approach loses interpretability of parameters. Depending on the application, the computation of sensitivities with respect to intuitive parameters may be a stringent requirement. This disadvantage is generally shared by neural network based models in engineering applications.

More generally, in some situations dimensionality reduction (even without interpretability of parameters) may be a desirable additional feature that is not directly achieved by our operator deep smoothing approach. Combining the VAE method [5] with our operator deep smoothing approach could lead to further promising potential applications of neural operators. [18] introduces *neural mappings*, which generalize neural operators to mixed infinite-/finite-dimensionality for input or output spaces. This motivates a discretization-invariant GNO-based encoder, fit to handle raw incoming market data, and a classical decoder to extend the operator deep smoothing approach to a VAE-like architecture.

Acknowledgements

Access to data was made possible through the Clarify Project funding by an Imperial College SME engagement grant jointly with Zeliade Systems, as well as the MSc in Mathematics and Finance, Imperial College London. AJ gratefully acknowledges financial support from the EPSRC grants EP/W032643/1 and EP/T032146/1. Computational resources and support were provided by the Imperial College Research Computing Service [14]. For the purpose of open access, the author(s) has applied a Creative Commons Attribution (CC BY) licence (where permitted by UKRI, ‘Open Government Licence’ or ‘Creative Commons Attribution No-derivatives (CC BY-ND) licence’ may be stated instead) to any Author Accepted Manuscript version arising.

References

- [1] Cboe GLOBAL MARKETS REPORTS TRADING VOLUME FOR DECEMBER AND FULL YEAR 2023. Technical report, 2023.
- [2] D. Ackerer, N. Tagasovska, and T. Vatter. Deep smoothing of the implied volatility surface. In *Advances in NeurIPS*, volume 33, pages 11552–11563, 2020.
- [3] A. Anandkumar, K. Azizzadenesheli, K. Bhattacharya, N. Kovachki, Z. Li, B. Liu, and A. Stuart. Neural operator: Graph kernel network for partial differential equations. *ICLR 2020 Workshop on Integration of Deep Neural Models and Differential Equations*.
- [4] F. Baschetti, G. Borretti, and P. Rossi. Deep calibration with random grids. *Quantitative Finance*, pages 1–23, 2024.
- [5] M. Bergeron, N. Fung, J. Hull, and Z. Poulos. Variational autoencoders: A hands-off approach to volatility. *arXiv:2102.03945*, 2021.
- [6] M. Chataigner, S. Crépey, and M. Dixon. Deep local volatility. *Risks*, 8(3):82, 2020.
- [7] J. Corbetta, P. Cohort, I. Laachir, and C. Martini. Robust calibration and arbitrage-free interpolation of SSVI slices. *Decisions in Economics and Finance*, 42(2):665–677, 2019.
- [8] K. Doherty, I. Almeida, and E. Popina. Options Are the Hottest Trade on Wall Street. *Bloomberg.com*, Oct. 2023.
- [9] M. Fukasawa. The normalizing transformation of the implied volatility smile. *Mathematical Finance*, 22(4):753–762, 2012.
- [10] J. Gatheral. A parsimonious arbitrage-free implied volatility parameterization with application to the valuation of volatility derivatives. *Presentation at Global Derivatives & Risk Management, Madrid*, 2004.
- [11] J. Gatheral and A. Jacquier. Arbitrage-free SVI volatility surfaces. *Quantitative Finance*, 14(1):59–71, 2014.
- [12] G. Guo, A. Jacquier, C. Martini, and L. Neufcourt. Generalized arbitrage-free SVI volatility surfaces. *SIAM Journal on Financial Mathematics*, 7(1):619–641, 2016.
- [13] Z. Hao, Z. Wang, H. Su, C. Ying, Y. Dong, S. Liu, Z. Cheng, J. Song, and J. Zhu. GNOT: A general neural operator transformer for operator learning. In *International Conference on Machine Learning*, pages 12556–12569. PMLR, 2023.
- [14] M. Harvey. Imperial College Research Computing Service (<http://doi.org/10.14469/hpc/2232>). 2017.
- [15] S. Hendriks and C. Martini. The extended SSVI volatility surface. *SSRN 2971502*, 2017.
- [16] A. Hernandez. Model calibration with neural networks. *SSRN:2812140*, 2016.
- [17] B. Horvath, A. Muguruza, and M. Tomas. Deep learning volatility: a deep neural network perspective on pricing and calibration in (rough) volatility models. *Quantitative Finance*, 21(1):11–27, 2021.
- [18] D. Z. Huang, N. H. Nelsen, and M. Trautner. An operator learning perspective on parameter-to-observable maps. *arXiv:2402.06031*, 2024.
- [19] P. Jäckel. Let’s be rational. *Wilmott*, 2015(75):40–53, 2015.
- [20] N. Kovachki, Z. Li, B. Liu, K. Azizzadenesheli, K. Bhattacharya, A. Stuart, and A. Anandkumar. Neural operator: Learning maps between function spaces with applications to PDEs. *Journal of Machine Learning Research*, 24(89):1–97, 2023.
- [21] S. Lanthaler, Z. Li, and A. M. Stuart. The nonlocal neural operator: Universal approximation. *arXiv:2304.13221*, 2023.
- [22] Z. Li, N. Kovachki, K. Azizzadenesheli, B. Liu, K. Bhattacharya, A. Stuart, and A. Anandkumar. Fourier neural operator for parametric partial differential equations. *arXiv:2010.08895*, 2020.
- [23] Z. Li, H. Zheng, N. Kovachki, D. Jin, H. Chen, B. Liu, K. Azizzadenesheli, and A. Anandkumar. Physics-informed neural operator for learning partial differential equations. *ACM/JMS Journal of Data Science*, 2021.
- [24] L. E. Lingsch, M. Y. Michelis, E. de Bezenac, S. M. Perera, R. K. Katzschmann, and S. Mishra. A structured matrix method for nonequispaced neural operators. *arXiv:2305.19663*, 2023.

- [25] S. Liu, A. Borovykh, L. A. Grzelak, and C. W. Oosterlee. A neural network-based framework for financial model calibration. *Journal of Mathematics in Industry*, 9(1):9, 2019.
- [26] V. Lucic. Normalizing volatility transforms and parameterization of volatility smile. *SSRN:3835233*, 2021.
- [27] C. Martini and A. Mingone. No arbitrage SVI. *SIAM Journal on Financial Mathematics*, 13(1):227–261, 2022.
- [28] C. Martini and A. Mingone. Refined analysis of the no-butterfly-arbitrage domain for SSVI slices. *Journal of Computational Finance*, 27(2), 2023.
- [29] A. Mingone. No arbitrage global parametrization for the eSSVI volatility surface. *Quantitative Finance*, 22(12):2205–2217, 2022.
- [30] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimeshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library, 2019.
- [31] M. Roper. Arbitrage free implied volatility surfaces. <https://www.maths.usyd.edu.au/u/pubs/publist/preprints/2010/roper-9.pdf>, 2010.
- [32] A. Tran, A. Mathews, L. Xie, and C. S. Ong. Factorized Fourier neural operators. *arXiv:2111.13802*, 2021.
- [33] L. Van Mieghem, A. Papapantoleon, and J. Papazoglou-Hennig. Machine learning for option pricing: an empirical investigation of network architectures. *arXiv:2307.07657*, 2023.
- [34] Y. Zheng. Machine Learning and Option Implied Information. *PhD Thesis, Imperial College London*, 2018.

A Neural Operators

We give a review of the kernel integral transform neural operator framework of [20], and expand in more detail on its graph neural operator.

Notation and Terms Let \mathcal{A} and \mathcal{U} be input and output space of an operator learning problem, as introduced in Section 3.1. Then, \mathcal{A} and \mathcal{U} are Banach spaces of functions $D \rightarrow \mathbb{R}^{c_{\text{in}}}$ and $D \rightarrow \mathbb{R}^{c_{\text{out}}}$, respectively, where D is a (bounded) domain in \mathbb{R}^d . The mathematical analysis of neural operators in [20] summarized hereafter, as well as the definition of universality and discretization-invariance in Section 3.1, makes use of the following terms and notations.

In the context of [20], a domain is a bounded and connected open set that is topologically regular (in the sense that it is the interior of its closure). A domain D is *Lipschitz* if its boundary locally is the graph of a Lipschitz continuous function defined on an open ball of \mathbb{R}^{d-1} . An *open ball* – for any metric space $\mathcal{X} = (\mathcal{X}, d)$ – is the set

$$\mathbb{B}_{\mathcal{X}}(x, \varepsilon) = \{y \in \mathcal{X} : d(y, x) < \varepsilon\}.$$

A *discrete refinement* of \mathcal{X} is a nested sequence (π_n) of discretizations of \mathcal{X} (finite subsets of \mathcal{X}), such that for every $\varepsilon > 0$ there is $N \in \mathbb{N}$ such that $\{\mathbb{B}(x, \varepsilon) : x \in \pi_N\}$ covers \mathcal{X} .

We consider the space $C(\mathcal{A}, \mathcal{U})$ of continuous operators between \mathcal{A} and \mathcal{U} . $C(\mathcal{A}, \mathcal{U})$ is topologized by *uniform convergence on compact sets*. With respect to this topology, a sequence $(F_n)_{n \in \mathbb{N}}$ in $C(\mathcal{A}, \mathcal{U})$ converges with limit $F \in C(\mathcal{A}, \mathcal{U})$, if, for every $\varepsilon > 0$ and every compact set K in \mathcal{A} , it holds

$$\lim_{n \rightarrow \infty} \|F_n - F\|_{\infty, K} = 0.$$

Here,

$$\|H\|_{\infty, K} = \sup_{a \in K} \|H(a)\|_{\mathcal{U}}, \quad H \in C(\mathcal{A}, \mathcal{U}).$$

It is well-known that this topology on $C(\mathcal{A}, \mathcal{U})$ is induced by the metric

$$\rho(F, G) = \sum_{n=0}^{\infty} \frac{\|G - F\|_{\infty, \mathbb{B}_{C(\mathcal{A}, \mathcal{U})}(0, n)}}{1 \vee \|G - F\|_{\infty, \mathbb{B}_{C(\mathcal{A}, \mathcal{U})}(0, n)}}, \quad F, G \in C(\mathcal{A}, \mathcal{U}).$$

Therefore, the notion of *density* in $C(\mathcal{A}, \mathcal{U})$, as used to define universality of neural operators in Section 3.1, is well-defined.

A.1 Kernel Integral Neural Operators

Kernel Integral Transform Neural Operators and Universality A kernel integral transform neural operator consists of the sequential application of:

1) A *lifting layer*

$$L_{\mathcal{P}} : [a : D \rightarrow \mathbb{R}^{c_{\text{in}}}] \mapsto [h_0 : D \rightarrow \mathbb{R}^{c_0}, h_0(x) = \mathcal{P}(a(x))],$$

given by the pointwise application of a function $\mathcal{P} : \mathbb{R}^{c_{\text{in}}} \rightarrow \mathbb{R}^{c_0}$.

2) The forward propagation through J neural operator layers L_0, \dots, L_{J-1} :

$$[h_0 : D \rightarrow \mathbb{R}^{c_0}] \xrightarrow{L_0} [h_1 : D \rightarrow \mathbb{R}^{c_1}] \xrightarrow{L_1} \dots \xrightarrow{L_{J-1}} [h_{J-1} : D \rightarrow \mathbb{R}^{c_J}];$$

each layer L_j operates as

$$h_{j+1}(y) = \sigma_j \left(W_j h_j(y) + \int_D \kappa_j(y, x) h_j(x) dx + b_j(y) \right), \quad y \in D, \quad (12)$$

where

- $W_j \in \mathbb{R}^{c_{j+1} \times c_j}$ is a weight matrix applied pointwise,
- $\kappa_j \in C(D \times D, \mathbb{R}^{c_{j+1} \times c_j})$ is a kernel function parametrizing the integral transform and subject to integrability conditions,
- The bias term b is itself a function from D to $\mathbb{R}^{c_{j+1}}$,

- σ_j is a classical neural network activation function.

3) A projection layer

$$L_{\mathcal{Q}}: [h_J: D \rightarrow \mathbb{R}^{c_J}] \mapsto [u: D \rightarrow \mathbb{R}^{c_{\text{out}}}, u(x) = \mathcal{Q}(h_J(x))],$$

given by the pointwise application of a function $\mathcal{Q}: \mathbb{R}^{c_J} \rightarrow \mathbb{R}^{c_{\text{out}}}$.

In practice, all components (lifting, kernel functions, projection) are implemented as classical feedforward neural networks (FFNs). This neural operator architecture is universal in the following sense.

Theorem A.1 (Universal Approximation; Theorem 11 of [20]). *Let D be a (bounded) Lipschitz domain. Assume:*

- $\mathcal{A} = W^{k_1, p_1}(D)$ for $k \in \mathbb{N}_{\geq 0}$ and $1 \leq p_1 < \infty$, or $\mathcal{A} = C(\overline{D})$.
- $\mathcal{U} = W^{k_2, p_2}(D)$ for $k \in \mathbb{N}_{\geq 0}$ and $1 \leq p_2 < \infty$, or $\mathcal{U} = C(\overline{D})$.

Then, a subset of kernel integral neural operators, with kernel functions and bias functions taken from a suitable set of FNNs, is dense in $C(\mathcal{A}, \mathcal{U})$.

Discretization-Invariant Implementations Consider a neural operator F^θ and let $\pi = \{x_l\}_{l=1}^p$ be a discretization of D . To make sense of a basic discretization-invariant implementation for F^θ , associate with π a partition (D_1, \dots, D_p) of D for which $\lambda_d(D_l) > 0$ and $x_l \in D_l$ for $l = 1, \dots, p$. Here λ_d denotes the Lebesgue measure on \mathbb{R}^d . Consider the following implementation of F^θ (written in terms of a single constituent layer $L = (W, \kappa, b, \sigma)$):

$$\tilde{L}(h|_\pi)(y) = \sigma \left(Wv(y) + \sum_{l=1}^p \kappa(y, x_l) h(x_l) \lambda_d(D_l) + b(y) \right), \quad y \in D. \quad (13)$$

[20] establishes the following.

Theorem A.2 (Discretization Invariance; Theorem 8 of [20]). *Let $F^\theta: \mathcal{A} \rightarrow \mathcal{U}$ be a kernel integral neural operator, where \mathcal{A} and \mathcal{U} both continuously embed into $C(\overline{D})$. Then, the implementation of F^θ based on (13) is discretization-invariant as defined in Section 3.1.*

(13) suggests the straightforward (quasi) Monte-Carlo inspired implementation

$$\tilde{L}(h|_\pi)(y) = \sigma \left(Wh(y) + \frac{\lambda_d(D)}{|\pi|} \sum_{x \in \pi} \kappa(y, x) h(x) + b(y) \right), \quad y \in D. \quad (14)$$

Most effectively, π is a low-discrepancy sequence in D .

A.2 Graph Neural Operator

The curse of dimensionality makes the direct implementation (14) prohibitively expensive in practice. Instead, [3] introduces *graph neural operators* (or, *GNOs*, for short) which replace the kernel integral operation at the heart of the framework by a sum approximation and organizes the constituent terms using a directed graph structure: The discretization $\pi = \{x_1, \dots, x_m\}$ of the input data \mathbf{a} is enriched with a directed graph structure $G_{\mathbf{a}} = (V, E)$, allowing the following implementation \tilde{F}^θ of F^θ :

$$\tilde{L}(v|_\pi)(y) = \sigma \left(Wv(y) + \frac{1}{|\mathcal{N}_{\text{in}}(y)|} \sum_{x \in \mathcal{N}_{\text{in}}(y)} \kappa(y, x) v(x) + b(y) \right), \quad y \in V. \quad (15)$$

Here, $\mathcal{N}_{\text{in}}(y)$ is the set of so-called *in-neighbors* of y in the graph $G_{\mathbf{a}}$: $x \in \mathcal{N}_{\text{in}}(y)$ iff $(x, y) \in E$. It is clear that, to compute output at y , the point y must be included as a node into the graph $G_{\mathbf{a}}$. It is important to reconcile the input and output locations of the layers when creating the graph structure to enable an efficient implementation using *message passing* algorithms.

The graph structure can be meticulously adjusted to implement various complexity-reducing techniques like *Nyström approximation* or *integration domain truncation* that effectively aim at a systematic reduction of the size of $\mathcal{N}_{\text{in}}(y)$; the naive implementation (13) is recovered for the case of a complete directed graph (with self-loops) for which $\mathcal{N}_{\text{in}}(y) = \pi$. Note that choosing $\mathcal{N}_{\text{in}}(y)$ as a strict subset of π breaks the guaranteed smoothness in the GNO output.

B Interpolation Graph Neural Operator

We detail our modifications of the GNO architecture. A generic Pytorch implementation is available at <https://github.com/rwic1/operator-deep-smoothing-for-implied-volatility>.

Let $\mathbf{v} = \{v(x)\}_{x \in \pi}$ be the given data.

Graph Construction Arguably part of the model architecture is the graph construction: Compile the set $\pi_{\text{out}} = \{y_l\}_{l=0}^q$ of points $y \in D$ at which to compute the smoothed surface $\hat{v}(y)$. During operator training, this will be the set of input locations (to compute the fitting loss) as well as any additional locations needed to compute auxiliary loss terms (the arbitrage losses $\mathcal{L}_{\text{but}}(\theta, \mathbf{v})$) and $\mathcal{L}_{\text{cal}}(\theta, \mathbf{v})$ in the case of volatility smoothing). For each $y \in \pi_{\text{out}}$, we construct the set of in-neighbors from the set of input data locations:

$$\mathcal{N}_{\text{in}}(y) \subseteq \pi_{\text{in}}. \quad (16)$$

In other words, we employ a Nyström approximation with nodes limited to the input data locations, which allows us to employ kernel functions with input skip connections. We set $G_{\mathbf{v}} = (\pi_{\text{out}}, E)$, where

$$E = \bigcup_{y \in \pi_{\text{out}}} \{(x, y) : x \in \mathcal{N}_{\text{in}}(y)\}.$$

Forward Propagation Given $G_{\mathbf{v}}$, we perform the first step of the forward propagation as follows:

$$\begin{cases} \tilde{h}_0(x) = \mathcal{P}_0(v(x)), & x \in \pi_{\text{in}} \\ h_1(y) = (\sigma_0 \circ \mathcal{Q}_0) \left(\mathcal{K}(\tilde{\mathbf{h}}_0; \mathbf{v})(y) + b_0 \right), & y \in \pi_{\text{out}}. \end{cases}$$

For the subsequent layers $j = 1, \dots, J - 1$, we then proceed using the classical scheme:

$$\begin{cases} \tilde{h}_j(y) = \mathcal{P}_j(h_j(y)), \\ h_{j+1}(y) = (\sigma_j \circ \mathcal{Q}_j) \left(W_j \tilde{h}_j(y) + \mathcal{K}_j(\tilde{\mathbf{h}}_j; \mathbf{v})(y) + b_j \right), & y \in \pi_{\text{out}}. \end{cases}$$

In the above:

- $\mathcal{P}_j : \mathbb{R}^{c_j} \rightarrow \mathbb{R}^{\tilde{c}_j}$ and $\mathcal{Q}_j : \mathbb{R}^{\tilde{c}_{j+1}} \rightarrow \mathbb{R}^{c_{j+1}}$ are layer-individual lifting and projection, implemented as FNNs.
- $W_j \in \mathbb{R}^{\tilde{c}_{j+1} \times \tilde{c}_j}$ is a weight matrix (not present for $j = 0$), while $b_j \in \mathbb{R}^{\tilde{c}_{j+1}}$ is a constant bias term.
- \mathcal{K}_j is the sum approximation of the kernel integral with kernel weight function $\kappa_j^W : D^2 \times \mathbb{R}^{\tilde{c}_j} \times \mathbb{R}^{c_0} \rightarrow \mathbb{R}^{\tilde{c}_{j+1} \times \tilde{c}_j}$ and kernel bias function $\kappa_j^b : D^2 \times \mathbb{R}^{\tilde{c}_j} \times \mathbb{R}^{c_0} \rightarrow \mathbb{R}^{\tilde{c}_{j+1}}$ (both with state and input skip connections):

$$\mathcal{K}_j(\tilde{\mathbf{h}}_j; \mathbf{v})(y) = \frac{1}{|\mathcal{N}_{\text{in}}(y)|} \sum_{x \in \mathcal{N}_{\text{in}}(y)} \kappa_j^W(y, x, \tilde{h}_j(x); v(x)) \tilde{h}_j(x) + \kappa_j^b(y, x, \tilde{h}_j(x); v(x)).$$

Both κ_j^W and κ_j^b are implemented as FNNs.

Note that omitting the local linear transformation in the first layer allows to extract the first hidden state $h_1(y)$ for all $y \in \pi_{\text{out}}$ from the lifted input $\tilde{\mathbf{h}}_0$, which is defined solely for the input locations $x \in \pi_{\text{in}}$. Providing each layer with its own lifting and projection allows to separate the hidden channel size c_0, \dots, c_J from the dimensions $\tilde{c}_0, \dots, \tilde{c}_J$ of the space in which the integral transformation is performed. Moreover, the individual lifting and projection help re-parametrize the state before performing the integral transform (inspired by the successful Transformer architecture), which allows to keep the size of the kernel weight matrix low.

C Supplementary Information: Data, Model, Training, Evaluation

This section contains additional information regarding our empirical study of the operator deep smoothing method for implied volatility smoothing.

C.1 Data

Data Source Our numerical experiments are based on the "Option Quotes" dataset product available for purchase from the CBOE. Our version of the dataset contains relevant data features for S&P 500 Index options for the years 2012 through 2021 and is summarized on a 20-minute interval basis.

Data Preparation We compute *simple mids* for options and underlying by averaging bid and ask quotes and use this aggregate as a reference price for all our subsequent computations. We calculate discount factors and forward prices from Put-Call parity using the industry standard technique based on linear regression. We compute time-to-expiry in units of one year as well as log-moneyness as defined in Section 2. We extract implied volatilities using the *py-vollib-vectorized* project available in the *Python Package Index* at the location <https://pypi.org/project/py-vollib-vectorized/>. *py-vollib-vectorized* implements a vectorized version of [19]’s *Let’s-be-rational* state-of-the-art method for computing implied volatility. We discard all implied volatilities of in-the-money options, or, in other words, we compose our implied volatility surface from Put options for non-positive log-moneyness values and Call options for positive log-moneyness values.

C.2 Model

We configure the modified GNO architecture introduced in Appendix B as follows:

- We impose the following restriction on the in-neighborhood sets $\mathcal{N}_{\text{in}}(y)$:

$$\mathcal{N}_{\text{in}}(y) \subseteq \{(\rho_l, z_l) \in \pi : |\rho_l - \rho| \leq \bar{\rho}\}. \quad (17)$$

This choice stems from the consideration that implied volatility smoothing requires limited global informational exchange along the time-to-expiry axis. We then generate $\mathcal{N}_{\text{in}}(y)$ of certain size K using a high-discrepancy subsampling heuristic detailed in the supplemented code material. Specifically, we use the values $\bar{\rho} = 0.3$ and $K = 50$.

- We employ three hidden layers and a channel size of 16: $J = 4$, and $c_1, c_2, c_3 = 16$ (c_0 and c_J are determined as 1 by the scalar dimension of volatility data). We use GELU-activations for the hidden layers and a Softplus-activation for the output layer (to ensure the positivity of the smoothed surfaces): $\sigma_0, \dots, \sigma_{J-1} = \text{GELU}$, and $\sigma_J = \text{Softplus}$.
- We retain $\mathcal{P}_0, \dots, \mathcal{P}_{J-1}$ and \mathcal{Q}_J as single-hidden layer FNNs with 64 hidden nodes and GELU-activations for the hidden layers. The remaining lifting and projections remain unutilized. In particular, $\tilde{c}_0, \dots, \tilde{c}_J = 16$.
- We implement the kernel weight and bias functions as two-hidden layer FNNs with 64 hidden nodes and GELU-activations for the hidden layers.

This configuration has been converged upon using manual experimentation and amounts to a total number of 102529 trainable parameters.

C.3 Training

Loss Function To ease notation we write $\hat{v}^\theta = \tilde{F}^\theta(\mathbf{v})$. We implement a Vega-weighted version of the fitting loss (9):

$$\mathcal{L}_{\text{fit}}(\theta; \mathbf{v}) = \left(\frac{1}{|\pi_{\mathbf{v}}|} \sum_{x \in \pi_{\mathbf{v}}} w_{\mathcal{V}}(x; \mathbf{v}) |(\hat{v}^\theta(x) - v(x))/v(x)|^2 \right)^{1/2}.$$

Here,

$$w_{\mathcal{V}}(x; \mathbf{v}) = \frac{\mathcal{V}(x, v(x))}{\frac{1}{|\pi|} \sum_{x \in \pi} \mathcal{V}(x, v(x))} \vee 1,$$

where $\mathcal{V}(x, v(x))$ is the Black-Scholes Vega defined in Section 1 (adjusted for transformed coordinates).

For the implementation of the no-arbitrage penalization terms \mathcal{L}_{but} and \mathcal{L}_{cal} , we first generate discretizations $\pi_\rho = \{\rho_1, \dots, \rho_m\}$ and $\pi_z = \{z_1, \dots, z_n\}$ of $[\rho_{\min}, \rho_{\max}]$ and $[z_{\min}, z_{\max}]$. We resolve the derivative terms $\partial_z v^\theta$ and $\partial_z^2 v^\theta$ on the synthetic rectilinear grid $\pi = \pi_\rho \times \pi_z$ using (central) finite differences. Then, we translate \mathcal{L}_{but} directly from (10) as

$$\mathcal{L}_{\text{but}}(\theta; \mathbf{v}, \pi) = \frac{1}{|\pi|} \sum_{x \in \pi} (\text{But}(x; \tilde{v}^\theta(x), \Delta_{z,\pi} \tilde{v}^\theta(x), \Delta_{z,\pi}^2 \tilde{v}^\theta(x)) - \varepsilon)^-,$$

where But is made consistent with the transformed coordinates and we used obvious notation for the finite differences. We use $\varepsilon = 10^{-3}$. On the other hand, we enforce the monotonicity constraint of Theorem 2.1 using

$$\mathcal{L}_{\text{cal}}(\theta; \mathbf{v}, \pi_\rho, \pi_z) = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \left(\frac{\tilde{v}^\theta(\rho_{i+1}, z_j)}{\tilde{v}^\theta(\rho_i, (\rho_{i+1} z_j) / \rho_i)} - \frac{\rho_i}{\rho_{i+1}} - \varepsilon \right)^-.$$

Compared to a derivative based implementation, this implementation is independent of the scale and – in our empirical experiments – has provided an improved signal. Since the Nyström approximation employed by the graph neural operator (as well as the choice (17)) break the guaranteed smoothness of the operator output, we additionally introduce $\|\partial_\rho^2 \hat{v}^\theta\|_2$ and $\|\partial_z^2 \hat{v}^\theta\|_2$ as regularization terms:

$$\mathcal{L}_{\text{reg-}\rho}(\theta; \mathbf{v}, \pi) = \sqrt{\frac{1}{|\pi|} \sum_{x \in \pi} |\Delta_{\rho,\pi}^2 \tilde{v}^\theta(x)|^2}, \quad \mathcal{L}_{\text{reg-}z}(\theta; \mathbf{v}, \pi) = \sqrt{\frac{1}{|\pi|} \sum_{x \in \pi} |\Delta_{z,\pi}^2 \tilde{v}^\theta(x)|^2}.$$

We compose the final loss function as a weighted sum of all terms introduced:

$$\mathcal{L}(\theta; \mathbf{v}, \pi_\rho, \pi_z) = \sum \begin{cases} \lambda_{\text{fit}} \mathcal{L}_{\text{fit}}(\theta; \mathbf{v}), \\ \lambda_{\text{but}} \mathcal{L}_{\text{but}}(\theta; \mathbf{v}, \pi_\rho \times \pi_z), \\ \lambda_{\text{cal}} \mathcal{L}_{\text{cal}}(\theta; \mathbf{v}, \pi_\rho, \pi_z), \\ \lambda_{\text{reg-}\rho} \mathcal{L}_{\text{reg-}\rho}(\theta; \mathbf{v}, \pi_\rho \times \pi_z), \\ \lambda_{\text{reg-}z} \mathcal{L}_{\text{reg-}z}(\theta; \mathbf{v}, \pi_\rho \times \pi_z). \end{cases}$$

The specific weights are

$$\begin{array}{ccccc} \lambda_{\text{fit}} & \lambda_{\text{cal}} & \lambda_{\text{but}} & \lambda_{\text{reg-}\rho} & \lambda_{\text{reg-}z} \\ \hline 1 & 10 & 10 & 0.01 & 0.01 \end{array}$$

The particular weighting of the individual terms has been retrieved by manual experimentation, led by the findings of [2].

Validation Loss Table 1 displays descriptive statistics of the validation losses.

	mean	std	min	25%	50%	75%	max
\mathcal{L}	0.062	0.032	0.029	0.053	0.059	0.064	0.506
\mathcal{L}_{fit}	0.022	0.016	0.009	0.016	0.020	0.024	0.232
\mathcal{L}_{but}	0.000	0.000	0.000	0.000	0.000	0.000	0.000
\mathcal{L}_{cal}	0.000	0.001	0.000	0.000	0.000	0.000	0.008
$\mathcal{L}_{\text{reg-}\rho}$	0.644	0.137	0.196	0.538	0.641	0.747	0.996
$\mathcal{L}_{\text{reg-}z}$	0.143	0.254	0.048	0.083	0.101	0.134	5.154

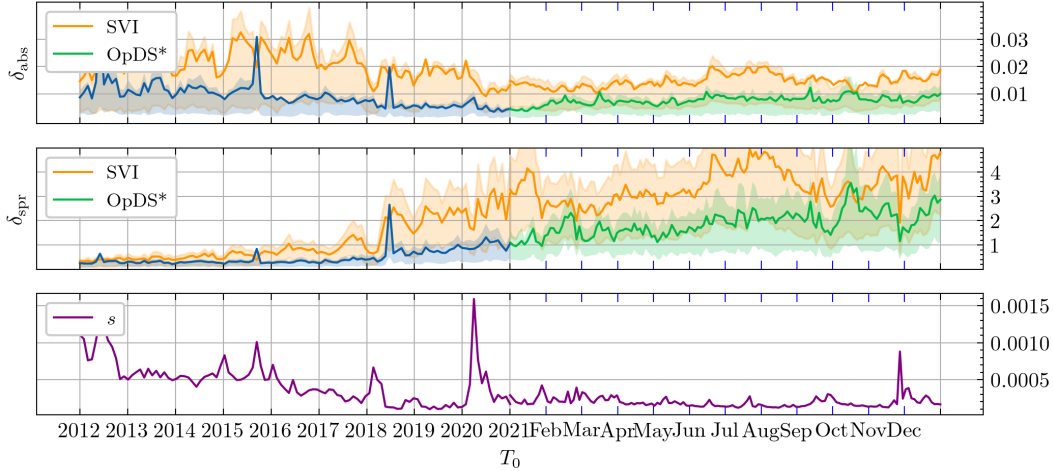


Figure 5: OpDS*: Benchmark metrics (surface-averages, between 25%- and 75%-quantiles) over training period (computed from validation dataset, resampled monthly) and testing period (resampled every two days).

C.4 Evaluation

Here we provide additional results to supplement our performance evaluation of operator deep smoothing (OpDS) for implied volatility. Table 2 and Table 3 display descriptive statistics of OpDS versus SVI. We also included a comparison to a less sophisticated version of OpDS without monthly finetuning, which we denote by OpDS*. Moreover, we include Figures 5 and 7, which are versions of Figures 3 and 4 with the only difference being that Figures 5 and 7 were created using OpDS* instead of OpDS. Finally, Figure 6 shows the average spatial distribution of benchmark metrics and arbitrage term over the training dataset, which complements the same averages on the test dataset shown in Figures 3.

Table 2: Descriptive statistics δ_{abs} over $\mathcal{D}_{\text{val}}/\mathcal{D}_{\text{test}}$.

	mean	std	1%	25%	50%	75%	99%
OpDS	0.009/0.005	0.007/0.001	0.003/0.003	0.006/0.004	0.008/0.005	0.010/0.005	0.021/0.007
OpDS*	0.009/0.007	0.007/0.001	0.003/0.003	0.006/0.07	0.008/0.007	0.010/0.008	0.021/0.012
SVI	0.021/0.015	0.006/0.002	0.007/0.010	0.016/0.013	0.020/0.014	0.025/0.016	0.034/0.020

Table 3: Descriptive statistics δ_{spr} over $\mathcal{D}_{\text{val}}/\mathcal{D}_{\text{test}}$.

	mean	std	1%	25%	50%	75%	99%
OpDS	0.479/1.265	0.662/0.347	0.193/0.609	0.274/1.025	0.330/1.240	0.550/1.451	1.526/2.453
OpDS*	0.479/1.866	0.662/0.574	0.193/0.731	0.274/1.457	0.330/1.826	0.550/2.233	1.526/3.445
SVI	1.124/3.382	0.877/0.826	0.301/1.464	0.492/2.827	0.715/3.320	1.646/3.914	3.710/5.247

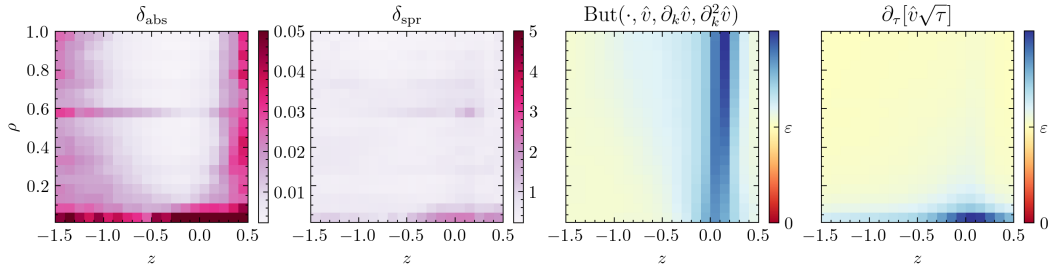


Figure 6: OpDS: Average spatial distribution of benchmark metrics and arbitrage term over train dataset.

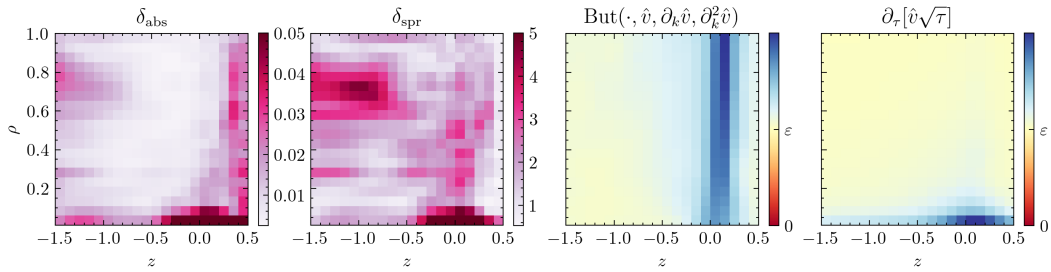


Figure 7: OpDS*: Average spatial distribution of benchmark metrics and arbitrage terms over test dataset (non-finetuned model).

D Additional Plots

D.1 Example Plots

We plot the results of operator deep smoothing (OpDS) vs. SVI on example inputs. To aid presentation, we include only every third market datapoint into the plots.

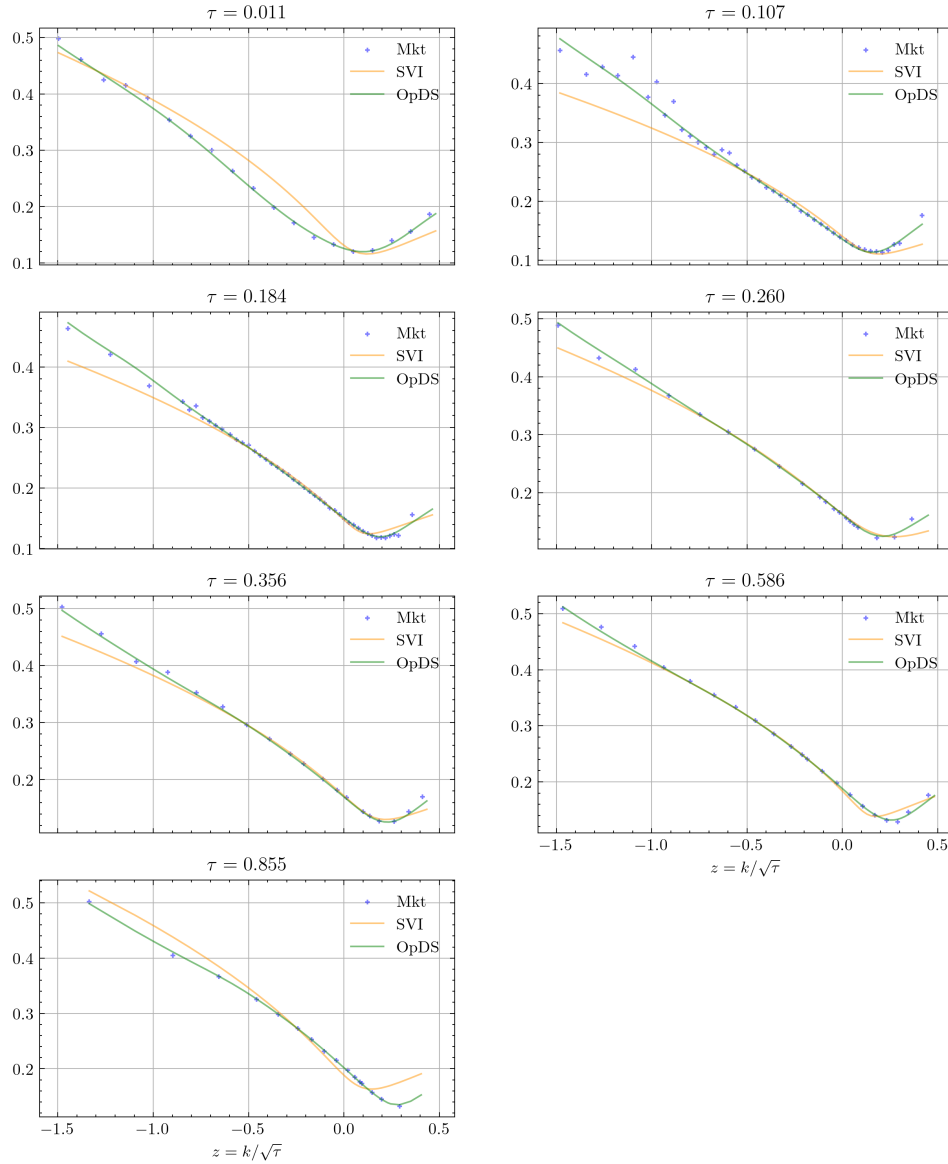


Figure 8: Smoothing of quotes v from 20.07.2012 at 10:50:00.

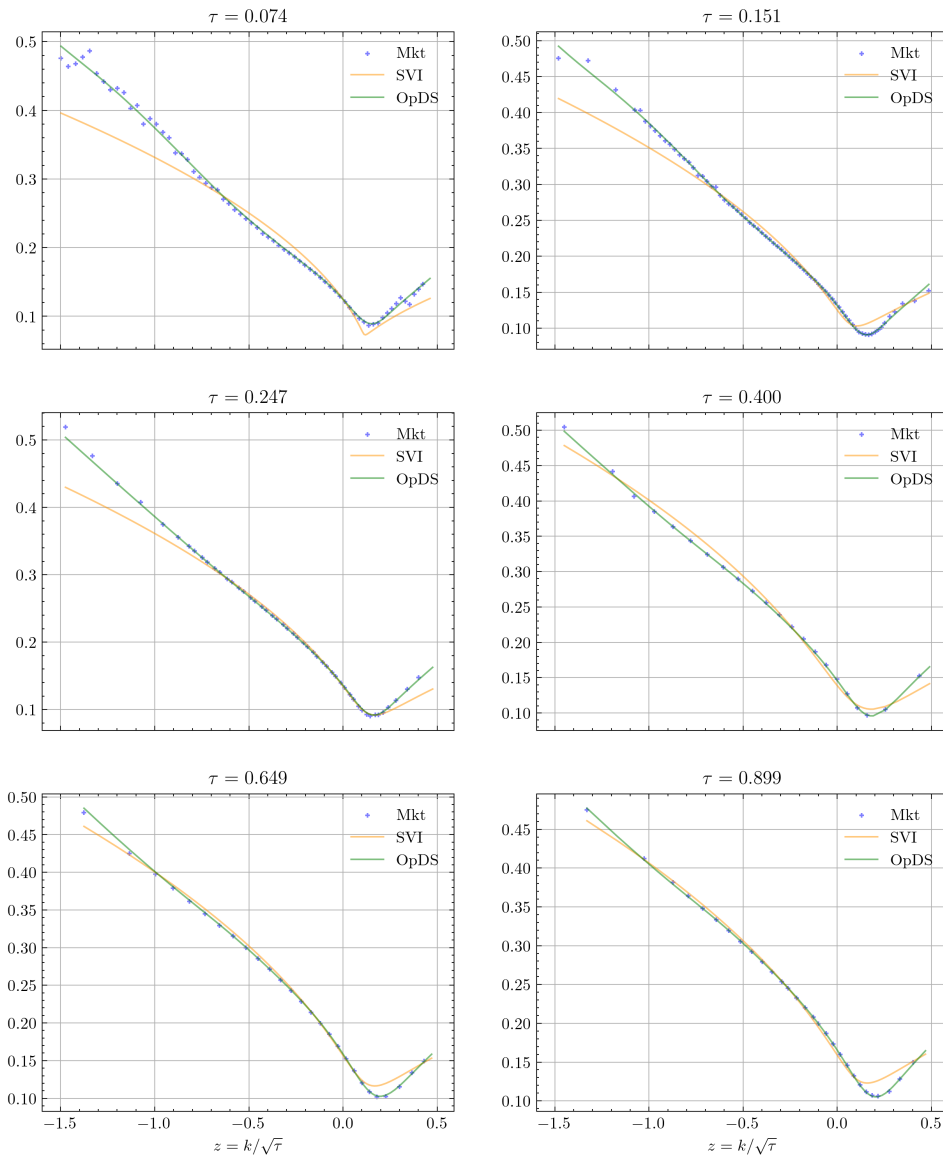


Figure 9: Smoothing of quotes v from 21.10.2016 at 13:10:00.

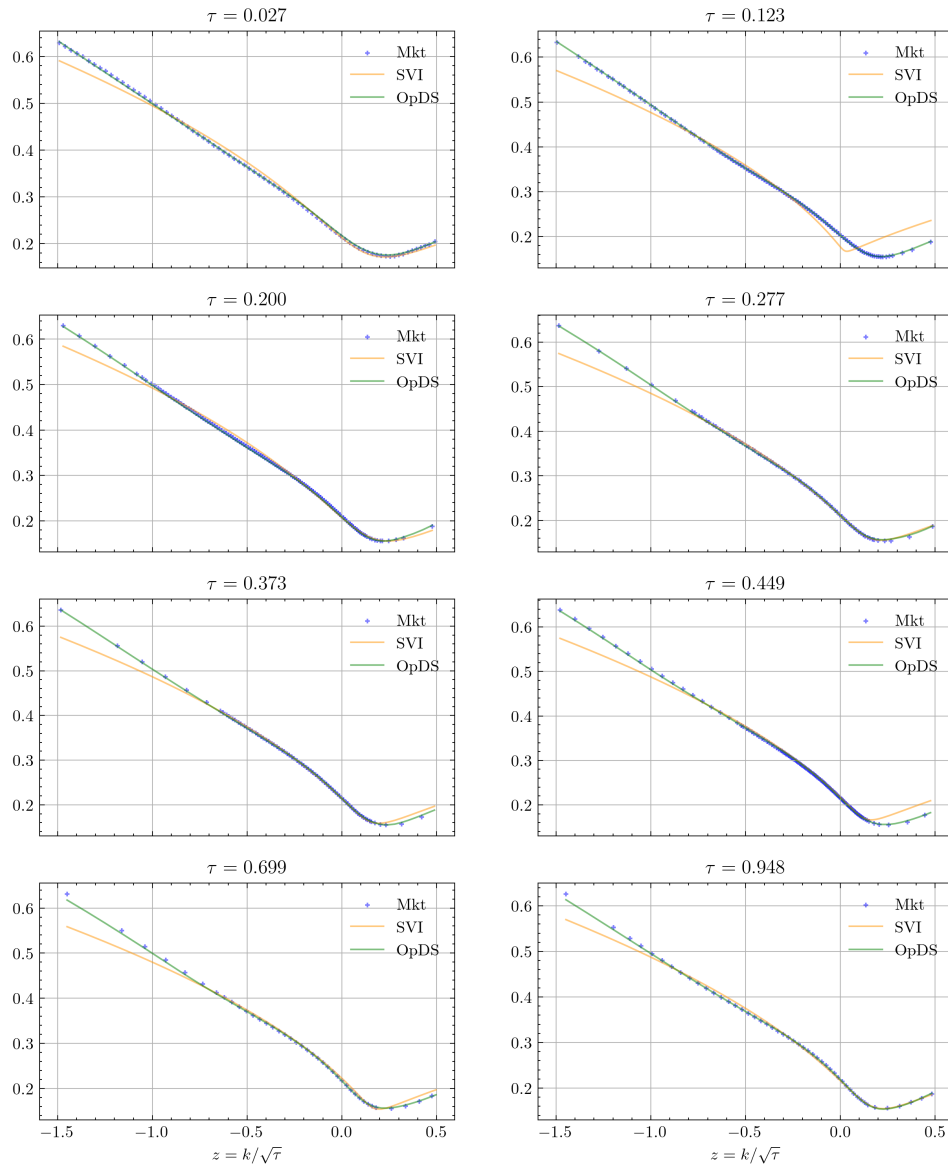


Figure 10: Smoothing of quotes v from 04.01.2021 at 10:50:00.

D.2 Monthly Breakdown of Spatial Distributions of Benchmark Metrics on Test Dataset

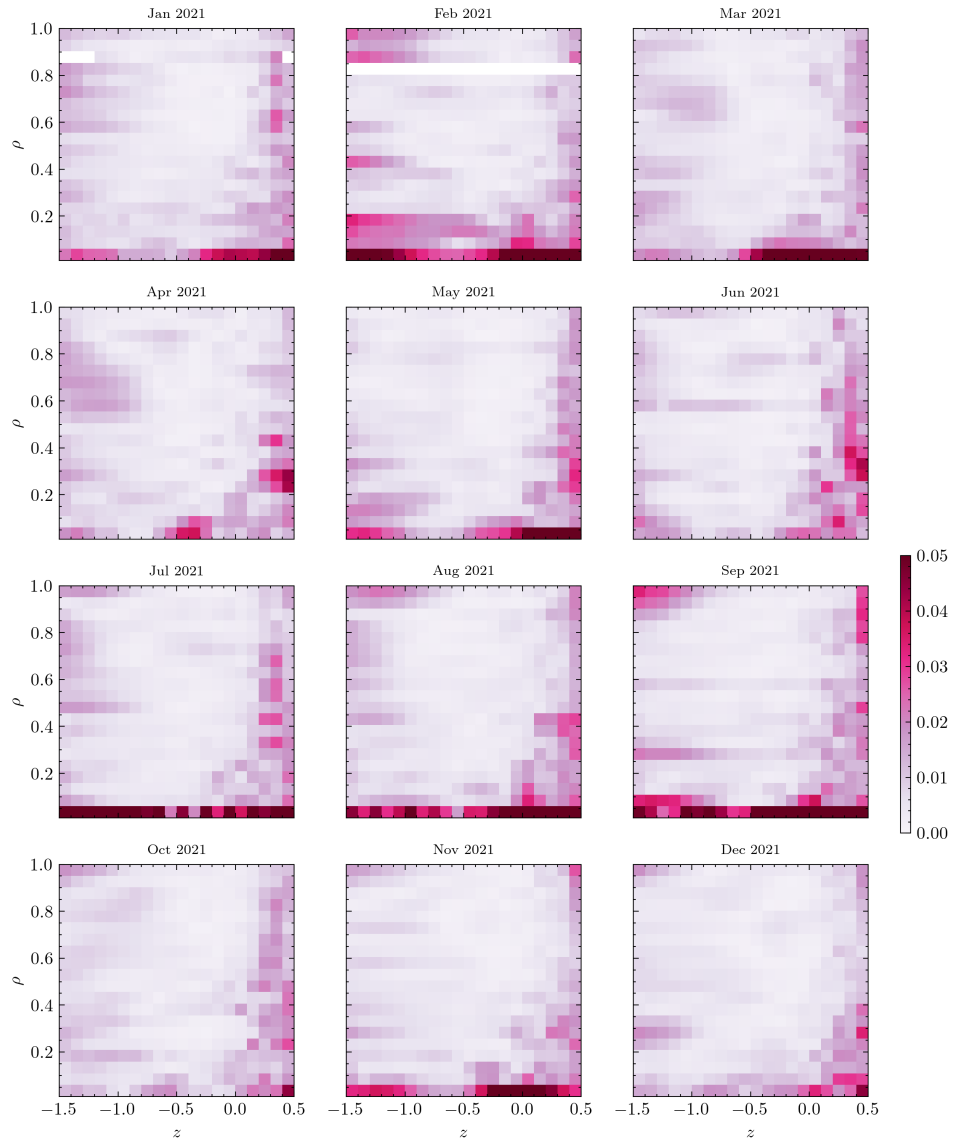


Figure 11: Average spatial distribution of δ_{abs} on $\mathcal{D}_{\text{test}}$ for OpDS with monthly finetuning, per month. Blank cells indicate that no data was available for the particular region in the respective month.

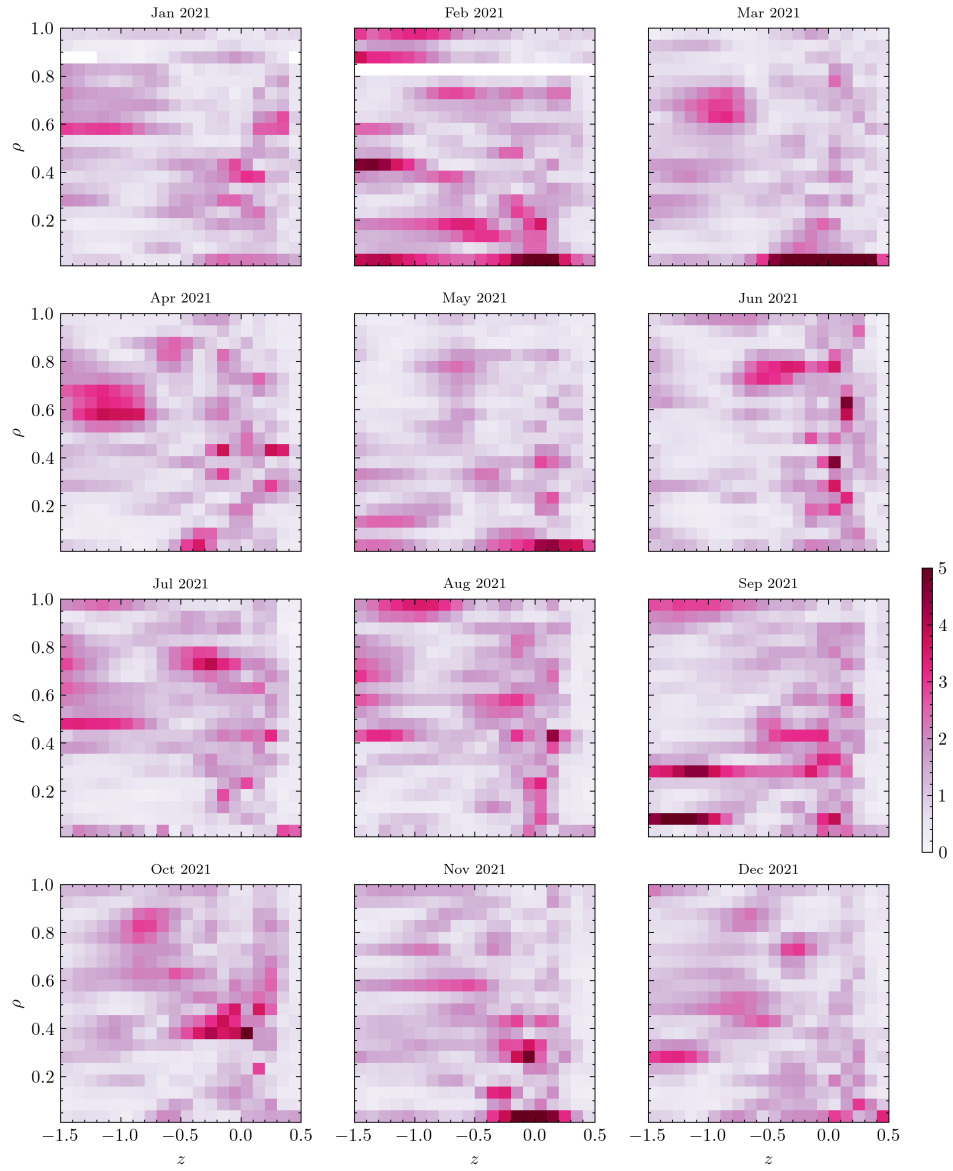


Figure 12: Average spatial distribution of δ_{fit} on D_{test} for OpDS with monthly finetuning, per month. Blank cells indicate that no data was available for the particular region in the respective month.

E SVI

SVI, originally devised by Gatheral [10], stands for *Stochastic Volatility Inspired* and is a low-dimensional parametrization for implied volatility slices (namely for each maturity). It captures the key features of implied volatility of Equity indices, and has become the industry-wide benchmark for implied volatility smoothing on such markets. Its "raw" variant parametrizes the "slice" of implied volatility at a given time-to-expiry τ as follows:

$$\hat{v}_\tau(k) = \sqrt{\frac{a + b \left(\rho(k - m) + \sqrt{(k - m)^2 + \sigma^2} \right)}{\tau}}, \quad \text{for all } k \in \mathbb{R},$$

where a, b, k, m, σ are real parameter values.

Calibration While stylistically accurate, SVI does not easily guarantee absence of static arbitrage opportunities, and several authors have investigated this issue [11, 27, 29, 28]. To produce our SVI benchmark we therefore rely on the constrained SLSQP optimizer provided by the *SciPy* scientific computing package for Python, with the mean square error objective, a positivity constraint and the constraint (3) (computed in closed form), and the following parameter bounds:

$$a \in \mathbb{R}, \quad b \in [0, 1], \quad \rho \in [-1, 1], \quad m \in [-1.5, 0.5], \quad \sigma \in [10^{-8}, 2].$$

We ignore the calendar arbitrage condition.