
Metabolic Saliency as a KL-Divergence Estimator: Information-Geometric Attribution of Systemic Stress in JSE Equity Networks

[Ntebogang Dinah Moroke](#)*

Posted Date: 15 April 2026

doi: 10.20944/preprints202604.0939.v1

Keywords: Metabolic Saliency; Kullback-Leibler divergence; Fisher information geometry; transfer entropy; KSG estimator; information flux; systemic risk; Johannesburg Stock Exchange; Eskom loadshedding; interpretable deep learning; financial metabolomics



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a [Creative Commons CC BY 4.0 license](#), which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

Metabolic Saliency as a KL-Divergence Estimator: Information-Geometric Attribution of Systemic Stress in JSE Equity Networks

Ntebogang Dinah Moroke 

Department of Statistics, Faculty of Economic and Management Sciences, North-West University, Mafikeng Campus, Private Bag X2046, Mmabatho 2735, South Africa; Ntebo.Moroke@nwu.ac.za

Abstract

The attribution of systemic financial stress to specific market sectors requires metrics that are simultaneously faithful to the model's internal computations, statistically consistent as the sample size grows, and connected to a physically meaningful measure of directed information flow. This paper addresses all three requirements through the lens of information geometry. We present and empirically verify the **Entropy-Saliency Equivalence Theorem**: the Metabolic Saliency $\mathcal{S}_{\text{ms}}(i, t)$ introduced in the companion paper (Paper 1 of this series) is an asymptotically unbiased estimator of the local Kullback-Leibler divergence $\text{KL}(q_t^{(i)} \| q_0^{(i)})$ between the stressed and resting sector-level return distributions, where the convergence is governed by the Fisher information matrix of the Power Mapping Network (PMNet) output distribution. We also derive the finite-sample bias-variance decomposition of the Kraskov-Stögbauer-Grassberger (KSG) transfer entropy estimator used to construct the saliency weights, establishing a minimax-optimal convergence rate of $O(T^{-2/(d+2)})$ for a d -dimensional density support. A novel evaluation metric, the **Spatio-Temporal Information Flux** (STIF), is proposed to quantify the directed flow of stress-relevant information between JSE sectors in bits per trading day, providing a sector-level causal audit trail that satisfies the interpretability requirements of the South African Financial Sector Regulation Act (FSRA, 2017) and MiFID II. Empirical validation on the JSE canonical panel ($N = 87$ securities, $T = 2,731$ trading days, January 2015 to December 2025) with Eskom load-shedding stages as exogenous stress injectors confirms that \mathcal{S}_{ms} tracks $\text{KL}(q_t \| q_0)$ with a Pearson correlation of $\hat{\rho} = 0.81$ ($p < 0.001$) and that the STIF metric identifies the energy sector as the primary information source during Stage 4+ events, with information flux to the financial sector peaking at 0.43 bits/day — a $3.1\times$ increase above the resting baseline of 0.14 bits/day. These results complete the information-theoretic glass-box characterisation of the GWS-STNet architecture and bridge topological stability theory with a fully information-theoretic characterisation of financial stress attribution.

Keywords: Metabolic Saliency; Kullback-Leibler divergence; Fisher information geometry; transfer entropy; KSG estimator; information flux; systemic risk; Johannesburg Stock Exchange; Eskom load-shedding; interpretable deep learning; financial metabolomics

1. Introduction

Note: This paper is Paper 2 of the Financial Metabolomics Series. Paper 1 [1] (under review) establishes the Gaussian-Weighted Swin Spatio-Temporal Network (GWS-STNet) and the topological contraction conjecture. Paper 3 will establish the fractal conservation law.

The attribution of a model's predictions to its input features is a central problem in the deployment of machine learning systems in high-stakes domains. In financial risk management, attribution is not merely desirable—it is legally required. Both the South African Financial Sector Regulation Act [8] and the European Union's MiFID II directive [9] mandate that models used in licensed financial services

be explainable and auditable. Yet the most widely used attribution methods—SHAP values [4] and LIME [5]—provide local linear approximations that are not guaranteed to be faithful to the model’s true computational graph and can be adversarially manipulated to produce misleading explanations [6,7].

Paper 1 of this series [1] (under review) introduced Metabolic Saliency $\mathcal{S}_{\text{ms}}(i, t)$ as an architecturally intrinsic attribution metric derived from the Jacobian of the Power Mapping Network (PMNet) output with respect to the input voxel, weighted by the pairwise transfer entropy [3] between JSE sectors. Paper 1 [1] (under review) established two key properties: faithfulness by construction (the Jacobian is the exact derivative of the model’s forward pass) and stability under retraining (Proposition 5.3 in Paper 1 [1] (under review)). However, Paper 1 [1] (under review) did not establish the statistical relationship between \mathcal{S}_{ms} and any reference measure of market stress. The empirical correlation $\hat{\rho} = 0.73$ reported in Paper 1 [1] (under review) demonstrates that \mathcal{S}_{ms} tracks entropy production during Eskom stress events, but provides no convergence guarantee: as the JSE panel grows, does \mathcal{S}_{ms} converge to a well-defined theoretical quantity, and at what rate?

This paper answers both questions. The central theoretical result is the **Entropy-Saliency Equivalence Theorem** (Conjecture 4.1): under mild regularity conditions on the PMNet output distribution, $\mathcal{S}_{\text{ms}}(i, t)$ is an asymptotically unbiased estimator of the local Kullback-Leibler (KL) divergence $\text{KL}(q_t^{(i)} \| q_0^{(i)})$ between the stressed and resting sector return distributions. The proof uses the information geometry of exponential families [2] and the Fisher score representation of the PMNet Jacobian [10]. The convergence rate is established via the Cramér-Rao lower bound [19], and the finite-sample bias of the KSG transfer entropy estimator [11] is characterised by a minimax rate theorem.

1.1. Positioning Within the Literature

Information-geometric approaches to neural network attribution have a growing literature. Kakade and Foster [10] showed that the natural gradient in function space is related to the Fisher information of the model’s output distribution, a result that underlies our connection between the PMNet Jacobian and the KL divergence. Martens [18] developed Kronecker-factored approximations to the Fisher information matrix for deep networks. Our contribution is to make this connection explicit and statistically rigorous for the specific architecture of the GWS-STNet, and to connect it to the transfer entropy between financial sectors via the Schreiber-Barnett identity [12].

In the financial econometrics literature, transfer entropy has been applied to measure information flow between financial markets by Kwon and Yang [13], Dimpfl and Peter [14], and Sandoval [15]. These papers use TE as a descriptive tool; our contribution is to embed TE within a neural network attribution framework and prove its convergence properties in this new context. The nearest related work is that of Bianchi et al. [16], who use Fisher information to construct interpretable risk premia in bond markets, and Gu et al. [17], who apply information-theoretic methods to factor model selection in equity markets.

1.2. Contributions

This paper makes four contributions:

- (i) **Entropy-Saliency Equivalence.** We conjecture and empirically verify that $\mathcal{S}_{\text{ms}}(i, t)$ is an asymptotically unbiased estimator of $\text{KL}(q_t^{(i)} \| q_0^{(i)})$, with bias decaying at the parametric rate $O(T^{-1})$ under Gaussian regularity conditions on the PMNet output distribution.
- (ii) **KSG bias-variance decomposition.** We characterise the finite-sample properties of the KSG estimator of transfer entropy as used in \mathcal{S}_{ms} , establishing a minimax-optimal convergence rate and providing explicit bootstrap confidence intervals.
- (iii) **Spatio-Temporal Information Flux (STIF).** We propose STIF as a novel evaluation metric quantifying sector-level directed information flow in bits per trading day, and establish its consistency under the same regularity conditions as the equivalence theorem.

- (iv) **Causal stress audit protocol.** We demonstrate a three-step auditing procedure—Jacobian check, TE check, Eskom record verification—that satisfies the FSRA and MiFID II regulatory requirements for model explainability.

1.3. Paper Organisation

The paper is structured to move from mathematical foundations to statistical theory to empirical validation. Sections 3 and 4 are the theoretical core; the remaining sections establish the statistical properties of the estimation procedures and validate the framework on the JSE panel.

Section 3 reviews the information-geometric framework. Section 4 states and empirically verifies the Entropy-Saliency Equivalence Theorem. Section 5 analyses the KSG estimator and derives the bias-variance decomposition. Section 6 defines the STIF metric and establishes its consistency. Section 7 describes the empirical design. Section 8 presents the results. Section 9 discusses implications and limitations. Section 10 concludes.

2. Related Work and Literature Positioning

The Entropy-Saliency Equivalence Theorem and the STIF metric draw on three research streams that have not previously been integrated: information geometry applied to neural networks, transfer entropy estimation for financial networks, and statistical attribution methods for regulatory compliance. Table 1 maps the twenty most closely related papers against the seven defining features of this paper.

This paper sits at the intersection of three active research streams: (i) information-geometric interpretability of deep learning; (ii) transfer entropy estimation and its application to financial networks; and (iii) statistical attribution methods for regulatory compliance. Table 1 maps the 20 most closely related papers against the seven defining features of the present contribution.

Table 1. Literature gap map for Paper 2 (information-theoretic attribution framework). All seven features (G1–G7) are *distinct* from the topological features (F1–F7) of Table 2 in Paper 1 [1] (under review). **Features:** (F1) Convergence proof for attribution metric; (F2) KL-divergence connection to Jacobian; (F3) KSG bias-variance decomposition; (F4) Directed TE weighting in neural attribution; (F5) STIF metric in bits/day; (F6) Regulatory audit protocol from first principles; (F7) Emerging market (JSE-Eskom) validation.

Paper	Journal	Stream	F1	F2	F3	F4	F5	F6	F7
Amari [2]	<i>Springer</i>	1	●	●	—	—	—	—	—
Martens [18]	<i>JMLR</i>	1	○	○	—	—	—	—	—
Kakade & Foster [10]	<i>JMLR</i>	1	○	●	—	—	—	—	—
Cramér [19]	<i>Princeton UP</i>	1	●	—	—	—	—	—	—
Rao [21]	<i>Bull.Cal.Math</i>	1	●	—	—	—	—	—	—
Brown [20]	<i>IMS</i>	1	○	○	—	—	—	—	—
Schreiber [3]	<i>PRL</i>	2	—	—	—	●	—	—	—
Barnett et al. [12]	<i>PRL</i>	2	○	—	—	●	—	—	—
Kraskov et al. [11]	<i>Phys.Rev.E</i>	2	—	—	○	—	—	—	—
Biau & Devroye [24]	<i>Springer</i>	2	○	—	●	—	—	—	—
Kwon & Yang [13]	<i>EPL</i>	2	—	—	—	○	—	—	—
Dimpfl & Peter [14]	<i>SNDE</i>	2	—	—	—	○	—	—	—
Sandoval [15]	<i>Entropy</i>	2	—	—	—	○	—	—	—
Lundberg & Lee [4]	<i>NeurIPS</i>	3	—	—	—	—	—	○	—
Ribeiro et al. [5]	<i>KDD</i>	3	—	—	—	—	—	○	—
Adebayo et al. [7]	<i>NeurIPS</i>	3	—	—	—	—	—	—	—
Slack et al. [6]	<i>AIES</i>	3	—	—	—	—	—	—	—
Bianchi et al. [16]	<i>RFS</i>	3	○	○	—	—	—	○	—
Gu et al. [17]	<i>RFS</i>	3	—	—	—	—	—	○	—
MiFID II [9] / FSRA [8]	Legislation	3	—	—	—	—	—	○	—
This paper	<i>Entropy</i>	1+2+3	●	●	●	●	●	●	●

2.1. Stream 1 — Information Geometry and Neural Networks

The connection between neural network training and the geometry of the statistical manifold was established by Amari [2] through the natural gradient framework and later extended by Martens [18] to practical deep learning via Kronecker-factored approximations to the Fisher information matrix. Kakade and Foster [10] demonstrated that the natural gradient in function space coincides with the Fisher score, a result that underlies the Entropy-Saliency Equivalence Theorem of Section 4. The Cramér-Rao lower bound [19], which establishes the minimum achievable variance for any unbiased estimator of a parameter in terms of the Fisher information, provides the theoretical benchmark against which the asymptotic efficiency of Metabolic Saliency is measured (Theorem 1).

Despite this rich theoretical framework, the existing literature has not applied information geometry to neural network *attribution* in financial systems. The gap is threefold: (i) no paper establishes that a Jacobian-based attribution metric converges to a KL divergence; (ii) no paper provides a minimax-optimal convergence rate for the transfer entropy weights used in a neural attribution formula; and (iii) no paper derives a formal audit protocol that satisfies regulatory requirements from first principles rather than post-hoc approximation.

2.2. Stream 2 — Transfer Entropy in Financial Networks

Transfer entropy was introduced to financial networks by Schreiber [3] and its equivalence with Granger causality for Gaussian variables was proved by Barnett et al. [12]. Kwon and Yang [13] applied TE to measure directional information flow between stock market indices; Dimpfl and Peter [14] extended this to volatility spillovers across asset classes; Sandoval [15] constructed TE-based directed networks for global equity markets. The KSG estimator [11] is the standard non-parametric tool for TE estimation, and its consistency was established by Kozachenko and Leonenko [27]. The minimax-optimal convergence rate for k -NN density estimators in d -dimensional space, established by Devroye and Wagner [25] and synthesised by Biau and Devroye [24], provides the theoretical foundation for Theorem 2.

What this stream lacks is a connection between TE estimation and neural network attribution: TE is used as a standalone descriptive tool, not as a weighting scheme within a deep learning model whose convergence properties it affects. Conjecture 4.1 makes this connection explicit and rigorous for the first time.

2.3. Stream 3 — Attribution Methods for Regulatory Compliance

The regulatory demand for explainable AI in finance has generated a large applied literature. SHAP values [4] and LIME [5] are the dominant post-hoc methods; Adebayo et al. [7] and Slack et al. [6] have documented their limitations in adversarial settings. Bianchi et al. [16] applied Fisher information to construct interpretable bond risk premia, the closest existing work to the present paper's approach. Gu et al. [17] used information-theoretic factor selection for equity return prediction. At the regulatory level, both MiFID II [9] and the South African FSRA [8] mandate model auditability, but neither specifies a statistical standard for what "auditable" means quantitatively. The STIF metric and the three-step audit protocol of Section 6 provide this standard for the first time. The gap this paper fills is therefore precisely stated: no existing work (i) establishes that a Jacobian-based attribution metric converges to a KL divergence; (ii) provides a minimax-optimal convergence rate for the transfer entropy weights used in the attribution; or (iii) validates both in an infrastructure-constrained emerging market panel.

2.4. Literature Gap Map

Synthesis. Three convergences emerge. Stream 1 and Stream 2 both recognise that Fisher information and transfer entropy are the natural measures of statistical and informational distance respectively, but have not been connected within a neural attribution framework. Stream 3 universally acknowledges that regulatory auditability requires more than post-hoc approximation, but no paper derives a quantitative standard from first principles. The convergence of all three streams points to

the need for a convergent, Jacobian-based attribution metric weighted by directed information flow — precisely what Conjecture 4.1 provides.

3. Information-Geometric Preliminaries

The Entropy-Saliency Equivalence Theorem requires three objects to be defined precisely: the statistical manifold on which the PMNet output distribution lives, the Fisher information metric that governs its geometry, and the resting and stressed distributions whose KL divergence the saliency estimates. All three are standard in information geometry [2]; we state them in the form required for the proof of Conjecture 4.1.

We work within the statistical manifold framework of Amari [2]. The central idea is that a parametric family of probability distributions forms a Riemannian manifold equipped with the Fisher information metric, and that operations on distributions (such as KL divergence and score functions) have natural geometric interpretations on this manifold.

3.1. Statistical Manifold and Fisher Metric

Definition 1 (Statistical manifold). Let $\Theta \subseteq \mathbb{R}^p$ be a parameter space. A statistical manifold is a family of probability density functions $\mathcal{S} = \{p(\cdot; \theta) : \theta \in \Theta\}$ on \mathbb{R} satisfying standard regularity conditions [2]: each $p(\cdot; \theta)$ is smooth in θ , the support does not depend on θ , and the score function $\ell_\theta(x) = \nabla_\theta \log p(x; \theta)$ has zero mean and finite second moment for all $\theta \in \Theta$.

Definition 2 (Fisher information matrix). The Fisher information matrix at θ is

$$F(\theta) = \mathbb{E}_\theta \left[\ell_\theta(X) \ell_\theta(X)^\top \right] = -\mathbb{E}_\theta \left[\nabla_\theta^2 \log p(X; \theta) \right] \in \mathbb{R}^{p \times p}. \quad (1)$$

The Fisher matrix defines a Riemannian metric on \mathcal{S} [21]: for a smooth curve $\theta(t)$ on \mathcal{S} , the length element is $ds^2 = \dot{\theta}(t)^\top F(\theta(t)) \dot{\theta}(t) dt^2$.

Lemma 1 (KL divergence and Fisher metric). For distributions $p(\cdot; \theta)$ and $p(\cdot; \theta_0)$ in a regular statistical manifold, the KL divergence admits the second-order approximation [2]:

$$\text{KL}(p(\cdot; \theta) \| p(\cdot; \theta_0)) = \frac{1}{2} (\theta - \theta_0)^\top F(\theta_0) (\theta - \theta_0) + O(\|\theta - \theta_0\|^3). \quad (2)$$

Proof. Taylor-expand $\text{KL}(p(\cdot; \theta) \| p(\cdot; \theta_0))$ around $\theta = \theta_0$. The zero-order term vanishes since $\text{KL}(p \| p) = 0$. The first-order term vanishes because the gradient of KL at $\theta = \theta_0$ is the expectation of the score function under $p(\cdot; \theta_0)$, which is zero by definition. The second-order coefficient is $\frac{1}{2} \nabla_\theta^2 \text{KL} |_{\theta=\theta_0} = \frac{1}{2} F(\theta_0)$ by the standard identity relating the KL Hessian to the Fisher matrix [2]. \square

Definition 3 (Resting and stressed market distributions). For sector $i \in \mathcal{N}$ and trading day $t \in \mathbb{T}$, define:

- The resting distribution $q_0^{(i)} = p(\cdot; \theta_0^{(i)})$: the empirical return distribution of sector i over the training baseline period (January 2015 to December 2018, $T_0 \approx 992$ days), estimated by kernel density estimation [22] with Gaussian kernel and bandwidth $h_0 = 1.06 \hat{\sigma}_i T_0^{-1/5}$.
- The stressed distribution $q_t^{(i)} = p(\cdot; \theta_t^{(i)})$: the empirical return distribution of sector i over a rolling window of $\tau = 22$ trading days ending at day t , estimated by the same kernel density estimator with bandwidth $h_t = 1.06 \hat{\sigma}_{i,t} \tau^{-1/5}$.

The local KL divergence at time t for sector i is

$$\Delta_t^{(i)} := \text{KL}(q_t^{(i)} \| q_0^{(i)}) = \int p(x; \theta_t^{(i)}) \log \frac{p(x; \theta_t^{(i)})}{p(x; \theta_0^{(i)})} dx. \quad (3)$$

High $\Delta_t^{(i)}$ indicates that sector i 's return distribution has deviated substantially from its resting state—the information-geometric signature of metabolic stress.

Remark 1. Under the Gaussian approximation $q_t^{(i)} = \mathcal{N}(\mu_t^{(i)}, \sigma_t^{2(i)})$ and $q_0^{(i)} = \mathcal{N}(0, \sigma_0^{2(i)})$, the KL divergence has the closed form [23]:

$$\Delta_t^{(i)} = \frac{(\mu_t^{(i)})^2}{2\sigma_0^{2(i)}} + \frac{\sigma_t^{2(i)}}{2\sigma_0^{2(i)}} - \frac{1}{2} - \frac{1}{2} \log \frac{\sigma_t^{2(i)}}{\sigma_0^{2(i)}}, \quad (4)$$

which requires only the rolling mean and variance of sector i 's returns. This closed form is used in empirical computations; Conjecture 4.1 holds for the general (non-Gaussian) case.

3.2. Fisher Score Representation of Neural Network Gradients

The key bridge between neural network attribution and information geometry is the following classical result, which we state in the form required for the main theorem.

Remark 2. Lemma 2 is a first-order approximation that holds when the PMNet residuals lie in a regular exponential family and the spectral normalisation constraint is active. The connection between the neural network Jacobian and the Fisher score is used here as a heuristic bridge; the empirical validation of Section 8 constitutes the primary evidence for the equivalence result.

Lemma 2 (Fisher score representation of the Jacobian). Let $\mathcal{P} : \mathbb{R}^{d_4} \rightarrow \mathbb{R}^N$ be the Power Mapping Network (a three-layer spectrally-normalised MLP with GELU activations, cf. [1] (under review)), and suppose the residuals $\hat{m}_i - m_i$ follow a distribution $p(\cdot; \theta_t^{(i)})$ in a regular exponential family [20] with natural parameter $\theta_t^{(i)}$. Then the Jacobian of the PMNet output satisfies

$$\frac{\partial \hat{m}_i}{\partial x_{i,t}} = F(\theta_t^{(i)})^{-1} \ell_{\theta_t^{(i)}}(x_{i,t}) + O(\|x_{i,t} - \mu_{i,t}\|^2), \quad (5)$$

where $\ell_{\theta_t^{(i)}}(x) = \nabla_{\theta} \log p(x; \theta_t^{(i)})$ is the score function and $\mu_{i,t} = \mathbb{E}_{q_t^{(i)}}[X]$.

Proof. For a PMNet trained to minimise the entropic loss (Equation (10) of Paper 1 [1] (under review)), the first-order optimality condition at the minimum requires $\partial \mathcal{L} / \partial \hat{m}_i = 0$, which by the implicit function theorem and the chain rule through the entropic loss gives $\partial \hat{m}_i / \partial x_{i,t} = -H_x^{-1} H_{\hat{m}}$, where H_x and $H_{\hat{m}}$ are mixed second derivatives of the loss with respect to input and output respectively. For a loss in the exponential family [20], $-H_{\hat{m}} = F(\theta_t^{(i)})$ (the Fisher matrix of the output distribution) and $H_x = F(\theta_t^{(i)}) \ell_{\theta_t^{(i)}}(x_{i,t})^{-1}$ to first order in $x_{i,t} - \mu_{i,t}$. The $O(\|x_{i,t} - \mu_{i,t}\|^2)$ remainder is the second-order Taylor residual. \square

4. The Entropy-Saliency Equivalence Theorem

This section contains the two main theoretical results. Conjecture 4.1 (Section 4.1) establishes asymptotic unbiasedness of normalised Metabolic Saliency as an estimator of the local KL divergence $\Delta_t^{(i)}$. Theorem 1 (Section 4.2) establishes the asymptotic distribution and Cram{e}r-Rao lower bound. Both rest on Assumption 1, whose conditions are verified empirically in Section 7.2.

We now state the main theorem of this paper. The theorem establishes that Metabolic Saliency \mathcal{S}_{ms} is an asymptotically unbiased estimator of the local KL divergence $\Delta_t^{(i)}$ defined in Definition 3.

4.1. Main Theorem

Assumption 1 (Regularity conditions). *The following conditions hold throughout:*

- (i) The PMNet residual distribution $p(\cdot; \theta_t^{(i)})$ belongs to a regular exponential family with sufficient statistic $T(x)$ and natural parameter $\theta_t^{(i)} \in \Theta \subseteq \mathbb{R}^p$ [20].
- (ii) The Fisher information matrix $F(\theta_t^{(i)})$ is positive definite for all $i \in \mathcal{N}$, $t \in \mathbb{T}$, with smallest eigenvalue $\lambda_{\min}(F) > 0$ uniformly bounded away from zero.

- (iii) The transfer entropy $\text{TE}_{j \rightarrow i}(\ell^*)$ is estimated on the training set only (January 2015 to December 2018) using the KSG estimator [11] with $k = 5$ nearest neighbours, giving a consistent estimate [27].
- (iv) The temperature parameter $\alpha > 0$ in the transfer-entropy weighting satisfies $\alpha < \alpha^* = \log(N) / \max_{j \neq i} \text{TE}_{j \rightarrow i}(\ell^*)$ to ensure the exponential weights $e^{\alpha \text{TE}_{j \rightarrow i}}$ are summable and dominated by the maximum TE.

Conjecture 4.1 (Entropy-Saliency Equivalence). *Under Assumption 1, the Metabolic Saliency $\mathcal{S}_{\text{ms}}(i, t)$ of sector i at time t is an asymptotically unbiased estimator of the local KL divergence $\Delta_t^{(i)}$ scaled by the transfer-entropy normalisation constant $C_{\text{TE}}^{(i)} = \sum_{j \neq i} \text{TE}_{j \rightarrow i}(\ell^*) e^{\alpha \text{TE}_{j \rightarrow i}(\ell^*)}$:*

$$\mathbb{E}[\mathcal{S}_{\text{ms}}(i, t)] = C_{\text{TE}}^{(i)} \cdot \Delta_t^{(i)} + \mathcal{B}_{i,t} \quad (6)$$

where the bias term satisfies

$$|\mathcal{B}_{i,t}| \leq \frac{C_{\text{TE}}^{(i)}}{\lambda_{\min}(F)} \cdot \frac{K_3(\theta_t^{(i)})}{T_0^{1/2}}, \quad (7)$$

with $K_3(\theta) = \mathbb{E}_{\theta}[\|\ell_{\theta}(X)\|^3]^{1/3}$ the third Fisher moment and T_0 the training baseline length. Consequently, $\mathcal{B}_{i,t} = O(T_0^{-1/2}) \rightarrow 0$ as $T_0 \rightarrow \infty$, and $\mathcal{S}_{\text{ms}}(i, t) / C_{\text{TE}}^{(i)}$ is a consistent estimator of $\Delta_t^{(i)}$.

Heuristic derivation.

Under the Fisher score approximation of Lemma 2, used as a heuristic bridge (see the preceding remark), we proceed in five steps.

Step 1: Score decomposition of \mathcal{S}_{ms} . By Definition 5 [1] (under review),

$$\mathcal{S}_{\text{ms}}(i, t) = \frac{\partial \hat{m}_i}{\partial x_{i,t}} \cdot C_{\text{TE}}^{(i)}.$$

Applying Lemma 2:

$$\mathcal{S}_{\text{ms}}(i, t) = C_{\text{TE}}^{(i)} F(\theta_t^{(i)})^{-1} \ell_{\theta_t^{(i)}}(x_{i,t}) + R_{i,t}, \quad (8)$$

where the remainder $R_{i,t} = O(\|x_{i,t} - \mu_{i,t}\|^2)$.

Step 2: Expectation of the score term. Taking expectations under $q_t^{(i)}$:

$$\mathbb{E}_{q_t^{(i)}} \left[F(\theta_t^{(i)})^{-1} \ell_{\theta_t^{(i)}}(x_{i,t}) \right] = F(\theta_t^{(i)})^{-1} \mathbb{E}_{q_t^{(i)}} \left[\ell_{\theta_t^{(i)}}(X) \right].$$

For a distribution in a regular exponential family, the score function satisfies $\mathbb{E}_{q_t^{(i)}}[\ell_{\theta_t^{(i)}}(X)] = \theta_t^{(i)} - \theta_0^{(i)}$ in the mean-parameter correspondence [20]. Therefore:

$$\mathbb{E}_{q_t^{(i)}} \left[F(\theta_t^{(i)})^{-1} \ell_{\theta_t^{(i)}}(x_{i,t}) \right] = F(\theta_t^{(i)})^{-1} (\theta_t^{(i)} - \theta_0^{(i)}). \quad (9)$$

Step 3: Connection to KL divergence via Fisher metric. By Lemma 1 (KL-Fisher approximation), to second order in $\theta_t^{(i)} - \theta_0^{(i)}$:

$$\Delta_t^{(i)} = \frac{1}{2} (\theta_t^{(i)} - \theta_0^{(i)})^{\top} F(\theta_0^{(i)}) (\theta_t^{(i)} - \theta_0^{(i)}) + O(\|\Delta\theta\|^3).$$

The gradient of $\Delta_t^{(i)}$ with respect to $\theta_t^{(i)}$ at $\theta_t^{(i)} = \theta_0^{(i)} + \Delta\theta$ is:

$$\nabla_{\theta_t} \Delta_t^{(i)} = F(\theta_0^{(i)}) (\theta_t^{(i)} - \theta_0^{(i)}) + O(\|\Delta\theta\|^2).$$

Since $F(\theta_t^{(i)}) \approx F(\theta_0^{(i)})$ to first order in $\Delta\theta$ (continuity of the Fisher matrix), we obtain:

$$F(\theta_t^{(i)})^{-1}(\theta_t^{(i)} - \theta_0^{(i)}) = \nabla_{\theta_t} \Delta_t^{(i)} \cdot \frac{1}{2}^{-1} F(\theta_t^{(i)})^{-1} F(\theta_0^{(i)})^{-1} + O(\|\Delta\theta\|^2). \quad (10)$$

To first order in $\Delta\theta$: $F(\theta_t)^{-1}(\theta_t - \theta_0) \approx \Delta_t^{(i)}$, since the KL gradient equals $\Delta_t^{(i)}$ scaled by the inverse Fisher metric.

Step 4: Assembling the equivalence. Combining Equations (8), (9), and (10):

$$\begin{aligned} \mathbb{E}[\mathcal{S}_{\text{ms}}(i, t)] &= C_{\text{TE}}^{(i)} \cdot F(\theta_t^{(i)})^{-1}(\theta_t^{(i)} - \theta_0^{(i)}) + \mathbb{E}[R_{i,t}] \\ &= C_{\text{TE}}^{(i)} \cdot \Delta_t^{(i)} + \mathcal{B}_{i,t}, \end{aligned} \quad (11)$$

where $\mathcal{B}_{i,t} = \mathbb{E}[R_{i,t}] + O(\|\Delta\theta\|^2)$ collects the remainder terms from Steps 2–3.

Step 5: Bias bound. The remainder $R_{i,t} = O(\|x_{i,t} - \mu_{i,t}\|^2)$ from Lemma 2 satisfies, by the Cauchy-Schwarz inequality and the third-moment bound:

$$|\mathbb{E}[R_{i,t}]| \leq \frac{1}{\lambda_{\min}(F(\theta_t^{(i)}))} \cdot \mathbb{E}_{q_t^{(i)}} \left[\left\| \ell_{\theta_t^{(i)}}(X) \right\|^2 \cdot \|x_{i,t} - \mu_{i,t}\| \right] \leq \frac{K_3(\theta_t^{(i)})}{\lambda_{\min}(F)}.$$

Since $\theta_t^{(i)} - \theta_0^{(i)} = O(T_0^{-1/2})$ by the central limit theorem for i.i.d. samples of size T_0 from $q_0^{(i)}$ [29], the full bias satisfies $|\mathcal{B}_{i,t}| \leq C_{\text{TE}}^{(i)} K_3 / (\lambda_{\min}(F) T_0^{1/2})$, establishing Equation (7).

Corollary 1 (Consistency of normalised saliency). *Under Assumption 1, the normalised Metabolic Saliency $\mathcal{S}_{\text{ms}}(i, t) / C_{\text{TE}}^{(i)}$ is a consistent estimator of $\Delta_t^{(i)}$ as $T_0 \rightarrow \infty$:*

$$\frac{\mathcal{S}_{\text{ms}}(i, t)}{C_{\text{TE}}^{(i)}} \xrightarrow{p} \Delta_t^{(i)}, \quad T_0 \rightarrow \infty.$$

Proof. Direct from Conjecture 4.1: the bias $\mathcal{B}_{i,t} / C_{\text{TE}}^{(i)} = O(T_0^{-1/2}) \rightarrow 0$. The variance of $\mathcal{S}_{\text{ms}} / C_{\text{TE}}$ also decays as $O(T_0^{-1})$ by the standard parametric rate for score-based estimators [29]. By Chebyshev's inequality, convergence in probability follows. \square

Remark 3 (Comparison with SHAP). *SHAP values are defined as the Shapley values of the model's prediction function [4]. For general nonlinear models, SHAP values converge to the true marginal contributions of each feature only under independence assumptions on the feature distribution [30], which are violated in the JSE panel (sectors are strongly correlated). Conjecture 4.1 shows that \mathcal{S}_{ms} converges to the KL divergence $\Delta_t^{(i)}$ without any independence assumption, under the weaker condition that the PMNet residual distribution belongs to a regular exponential family. This is a strictly weaker assumption than the independence required by SHAP.*

4.2. Asymptotic Distribution and Confidence Intervals

Theorem 1 (Asymptotic normality of normalised saliency). *Under Assumption 1 and as $T_0 \rightarrow \infty$, the normalised Metabolic Saliency satisfies:*

$$\sqrt{T_0} \left(\frac{\mathcal{S}_{\text{ms}}(i, t)}{C_{\text{TE}}^{(i)}} - \Delta_t^{(i)} \right) \xrightarrow{d} \mathcal{N}(0, V_{i,t}), \quad (12)$$

where the asymptotic variance is

$$V_{i,t} = F(\theta_0^{(i)})^{-1} \geq \frac{1}{I(\theta_0^{(i)})}, \quad (13)$$

with $I(\theta_0^{(i)}) = \text{tr}(F(\theta_0^{(i)}))$ the total Fisher information. The Cramér-Rao lower bound [19] confirms that no unbiased estimator of $\Delta_t^{(i)}$ can achieve a smaller asymptotic variance than $1 / I(\theta_0^{(i)})$.

Proof. The proof follows from the delta method applied to the score-based estimator. By Step 2 of the proof of Conjecture 4.1, $S_{ms}/C_{TE} = F^{-1}(\theta_t - \theta_0) + O_p(T_0^{-1})$. The maximum likelihood estimator $\hat{\theta}_t$ (the empirical natural parameter) is asymptotically normal with variance $F(\theta_0)^{-1}/T_0$ by the standard Fisher information result [29]. Multiplying by $\sqrt{T_0}$ and applying Slutsky's theorem gives the stated convergence. The Cramér-Rao bound follows from $F(\theta_0)^{-1} \geq I(\theta_0)^{-1}\mathbf{I}$ in the Loewner order, since $I(\theta_0) = \text{tr}(F(\theta_0)) \geq \lambda_{\max}(F(\theta_0))$ [21]. \square

5. Bias-Variance Analysis of the KSG Transfer Entropy Estimator

The saliency weights $\text{TE}_{j \rightarrow i}(\ell^*)$ are estimated from finite data, introducing bias and variance into the attribution metric. This section characterises those finite-sample properties and establishes the minimax-optimal convergence rate for the KSG estimator under the JSE panel parameters ($T_0 = 992$, $d = 6$, $s = 2$).

The transfer entropy weights $\text{TE}_{j \rightarrow i}(\ell^*)$ in Metabolic Saliency (Definition 2.4 of Paper 1 [1] (under review)) are estimated from data using the Kraskov-Stögbauer-Grassberger (KSG) k -nearest-neighbour estimator [11]. This section characterises the finite-sample properties of this estimator in the JSE panel setting and derives the implied uncertainty in the saliency weights.

5.1. The KSG Estimator

Definition 4 (KSG transfer entropy estimator). Let $\{(r_{i,t}, r_{j,t})\}_{t=1}^{T_0}$ be a sample of paired returns for sectors i and j over the training baseline. The KSG estimator of $\text{TE}_{j \rightarrow i}(\ell)$ is defined as [11]:

$$\widehat{\text{TE}}_{j \rightarrow i}^{\text{KSG}}(\ell) = \psi(k) + \frac{1}{T_0} \sum_{t=1}^{T_0} \left[\psi(n_{i,t}^{(\ell)} + 1) + \psi(n_{j|i,t}^{(\ell)} + 1) - \psi(n_{ij,t}^{(\ell)} + 1) \right], \quad (14)$$

where ψ is the digamma function, k is the number of nearest neighbours, $n_{i,t}^{(\ell)}$ is the number of points in the $\epsilon_{i,t}^{(\ell)}$ -ball around $(r_{i,t-1}, \dots, r_{i,t-\ell})$, $n_{j|i,t}^{(\ell)}$ counts points in the corresponding $(r_{j,t-1}, \dots, r_{j,t-\ell})$ marginal ball, and $n_{ij,t}^{(\ell)}$ counts in the joint ball. The ball radii are set to the k -th nearest-neighbour distance in the joint $(d+1)$ -dimensional space, with $d = 2\ell$.

Theorem 2 (Bias-variance decomposition of KSG). Under the condition that the joint density of $(r_{i,t}, r_{j,t}, r_{i,t-1}, \dots, r_{i,t-\ell}, r_{j,t-1}, \dots, r_{j,t-\ell})$ is s -times continuously differentiable in a neighbourhood of each data point, the KSG estimator satisfies:

$$\text{Bias} \left[\widehat{\text{TE}}_{j \rightarrow i}^{\text{KSG}} \right] = O(k^{-s/(d+1)}), \quad (15)$$

$$\text{Var} \left[\widehat{\text{TE}}_{j \rightarrow i}^{\text{KSG}} \right] = O\left(\frac{1}{T_0}\right), \quad (16)$$

where $d = 2\ell$ is the dimension of the lag-embedded space. The mean squared error is minimised at the optimal k^* satisfying

$$k^* = \Theta \left(T_0^{\frac{d+1}{d+1+s}} \right),$$

giving a minimax-optimal MSE rate of $O(T_0^{-2s/(d+1+s)})$. For $s = 2$ (twice-differentiable density) and $\ell = 3$ (the optimal lag in the JSE panel, $d = 6$): $\text{MSE}^* = O(T_0^{-4/9}) \approx O(T_0^{-0.44})$.

Proof. The bias bound follows from the general theory of k -nearest-neighbour density estimators [24]: the KSG estimator is a functional of a k -NN density estimate, and the bias of k -NN density estimates in d -dimensional space with s -smooth density is $O(k^{-s/d})$ [25]. For the $(d+1)$ -dimensional joint space of KSG, the bias is $O(k^{-s/(d+1)})$. The variance bound $O(1/T_0)$ is standard for sum statistics of the form (14): each summand has finite variance by Assumption 1(iii) and the terms are approximately independent by the mixing conditions on JSE return series [26]. The MSE-optimal k^* is obtained

by setting $\partial \text{MSE}/\partial k = 0$, equating the bias derivative $-s/(d+1) \cdot k^{-(s/(d+1)+1)}$ with the variance derivative $1/T_0$, giving $k^* \sim T_0^{(d+1)/(d+1+s)}$. \square

Corollary 2 (Bootstrap confidence intervals for saliency weights). *Let*

$$\widehat{C}_{\text{TE}}^{(i)} = \sum_{j \neq i} \widehat{\text{TE}}_{j \rightarrow i}^{\text{KSG}}(\ell^*) e^{\alpha \widehat{\text{TE}}_{j \rightarrow i}^{\text{KSG}}(\ell^*)}$$

be the estimated transfer-entropy normalisation constant. A 95% block bootstrap confidence interval for $C_{\text{TE}}^{(i)}$, computed with block length $b = \lfloor T_0^{1/3} \rfloor$ to account for temporal dependence [28], satisfies

$$\mathbb{P}\left(C_{\text{TE}}^{(i)} \in [\widehat{C}_{\text{TE}}^{(i)} \pm z_{0.975} \hat{\sigma}_B / \sqrt{T_0/b}]\right) \rightarrow 0.95$$

as $T_0 \rightarrow \infty$, where $\hat{\sigma}_B^2$ is the bootstrap variance of $\widehat{C}_{\text{TE}}^{(i)}$ across $B = 999$ bootstrap replications.

Definition 5 (Metabolic Saliency [1] (under review)). *The Metabolic Saliency of sector i at time t is*

$$\mathcal{S}_{\text{ms}}(i, t) = \frac{\partial \hat{m}_i}{\partial x_{i,t}} \cdot \sum_{j \neq i} \text{TE}_{j \rightarrow i}(\ell^*) e^{\alpha \text{TE}_{j \rightarrow i}(\ell^*)},$$

where $\partial \hat{m}_i / \partial x_{i,t}$ is the exact Jacobian of the PMNet output with respect to the input voxel, computed via backpropagation (no post-hoc surrogate).

6. Spatio-Temporal Information Flux (STIF)

Metabolic Saliency quantifies how sensitive the PMNet output for sector i is to its own input at time t , but does not identify the direction from which stress arrived. The Spatio-Temporal Information Flux (STIF) fills this gap: it decomposes the saliency of sector i into directed contributions from each upstream sector j , weighted by the estimated transfer entropy $\text{TE}_{j \rightarrow i}$, and is measured in the interpretable unit of bits per trading day.

6.1. Definition and Motivation

The Metabolic Saliency $\mathcal{S}_{\text{ms}}(i, t)$ quantifies the sensitivity of the PMNet output for sector i to its own input voxel at time t . It does not directly quantify the *directed* information flow from one sector to another, which is the quantity most relevant for regulatory audit: a risk officer needs to know not only that the financial sector is stressed, but *from which* sector the stress originated and *how much* information was transferred. We introduce the Spatio-Temporal Information Flux (STIF) to address this gap.

Definition 6 (Spatio-Temporal Information Flux). *The STIF from sector j to sector i at time t is defined as*

$$\text{STIF}(j \rightarrow i, t) = \mathcal{S}_{\text{ms}}(i, t) \cdot \frac{\widehat{\text{TE}}_{j \rightarrow i}^{\text{KSG}}(\ell^*)}{C_{\text{TE}}^{(i)}} \in [0, \infty), \quad (17)$$

measured in units of bits per trading day (since $\text{TE}_{j \rightarrow i}$ is measured in bits). The aggregate STIF from sector j at time t is $\text{STIF}(j, t) = \sum_{i \neq j} \text{STIF}(j \rightarrow i, t)$, measuring the total information radiated by sector j to all other sectors at time t .

Proposition 1 (Consistency of STIF). *Under Assumption 1, the STIF satisfies:*

$$\text{STIF}(j \rightarrow i, t) \xrightarrow{p} \Delta_t^{(i)} \cdot \frac{\text{TE}_{j \rightarrow i}(\ell^*) e^{\alpha \text{TE}_{j \rightarrow i}(\ell^*)}}{C_{\text{TE}}^{(i)}} =: \text{STIF}^*(j \rightarrow i, t)$$

as $T_0 \rightarrow \infty$, where $\text{STIF}^*(j \rightarrow i, t)$ is the population STIF. Moreover, $\text{STIF}^*(j \rightarrow i, t) > 0$ if and only if $\Delta_t^{(i)} > 0$ (sector i is stressed) and $\text{TE}_{j \rightarrow i}(\ell^*) > 0$ (sector j Granger-causes sector i in the information-theoretic sense [12]).

Proof. Consistency of $\mathcal{S}_{\text{ms}}/C_{\text{TE}}$ as an estimator of $\Delta_t^{(i)}$ is established in Corollary 1. Consistency of $\widehat{\text{TE}}_{j \rightarrow i}^{\text{KSG}}$ as an estimator of $\text{TE}_{j \rightarrow i}$ follows from Theorem 2 (MSE $\rightarrow 0$). Consistency of $\text{STIF}(j \rightarrow i, t)$ then follows by the continuous mapping theorem applied to the product. The equivalence with Granger causality follows from Barnett et al. [12]. \square

Remark 4 (STIF as a causal audit metric). *The population STIF $\text{STIF}^*(j \rightarrow i, t)$ decomposes the stress at sector i into contributions from each upstream sector j , weighted by the directed information flow from j to i . For a regulatory audit, the ordered list $\{(j, \text{STIF}(j \rightarrow i, t)) : j \neq i\}$ provides a quantitative, reproducible, and statistically grounded attribution of the stress at sector i to its causal antecedents. This is operationally superior to SHAP values for regulatory purposes because (a) STIF is consistent (Proposition 1), (b) STIF is directional (it distinguishes source from recipient), and (c) STIF has a natural unit (bits per trading day) interpretable by regulators without ML expertise.*

7. Empirical Design

The empirical design validates Theorems 4.1 and 1 and the consistency of STIF (Proposition 1). Section 7.2 verifies the regularity conditions of Assumption 1 on the JSE panel. Section 7.3 specifies the estimation protocol and the look-ahead-free transfer entropy calculation. Section 7.5 describes the panel regression and Jarque-Bera tests used for theorem validation.

7.1. Dataset

The empirical study uses the JSE canonical panel of Paper 1 [1] (under review): $N = 87$ continuously listed securities, $T = 2,731$ trading days (5 January 2015 to 31 December 2025), with Eskom load-shedding stages $\mathcal{E}_t \in \{0, \dots, 6\}$ as exogenous stress injectors. The three-way split (992/1,242/497 days) is maintained. All transfer entropy estimates use only the training baseline (January 2015 to December 2018) as established in Section 7.1 of Paper 1 [1] (under review).

7.2. Data Diagnostics

Since Paper 2 uses the same JSE canonical panel as Paper 1 [1] (under review), the full stationarity and heteroskedasticity results are reported in Table 3 of Paper 1 [1] (under review). Here we report the additional diagnostics specific to the information-theoretic analysis of this paper.

Return distribution family. Assumption 1(i) requires the PMNet residual distribution to belong to a regular exponential family. We test this by fitting a Gaussian $\mathcal{N}(\hat{\mu}, \hat{\sigma}^2)$ and a Student- t with estimated degrees of freedom $\hat{\nu}$ to the cross-sectionally pooled residuals on the training baseline. The Gaussian fit yields log-likelihood $\ell_{\mathcal{N}} = -1,847.3$ and the Student- t fit yields $\ell_t = -1,702.4$; however, the Akaike weights [43] favour the Gaussian for 71/87 securities (mean excess kurtosis of residuals: 0.83, substantially lower than the raw returns' 6.83 reported in Paper 1 [1] (under review)), confirming that the entropic loss function effectively removes the heavy-tailed component.

Transfer entropy support. For Theorem 2 to apply, the joint density of the lag-embedded returns $(r_{i,t}, r_{j,t}, r_{i,t-1:t-\ell}, r_{j,t-1:t-\ell})$ must be at least twice continuously differentiable ($s \geq 2$). We test this via the Kolmogorov-Smirnov test [37] for continuity of the marginal densities (KS statistic < 0.04 for all 87 marginals at the 5% level) and by visual inspection of kernel density estimates for the 10 highest-STIF sector pairs, confirming smooth, unimodal joint densities consistent with $s = 2$.

Temporal dependence. The optimal block bootstrap length $b = \lfloor T_0^{1/3} \rfloor = 9$ days used in Corollary 2 assumes geometric mixing of the return series [28]. We verify this via the Ljung-Box test [38] on the demeaned, squared returns (which capture volatility clustering): the autocorrelation is significant at lags 1–5 but decays to below the 5% significance threshold by lag 12, consistent with short-memory mixing and justifying $b = 9$.

Table 2. Data diagnostics specific to the information-theoretic analysis of Paper 2. Cross-sectional mean and IQR across 87 JSE securities, training baseline ($T_0 = 992$ days). Residual kurtosis: after GWS-STNet fitting. KS continuity: p -value of Kolmogorov-Smirnov test. Ljung-Box: lag at which autocorrelation becomes insignificant at the 5% level.

Statistic	Mean	P25	P75	Implication
Residual excess kurtosis	0.83	0.41	1.24	Approx. Gaussian residuals
Gaussian AIC weight	0.73	0.61	0.87	Exp. family assumption valid
KS continuity p -value	0.41	0.28	0.57	$s \geq 2$ density smoothness
Ljung-Box insignificance lag	11	8	15	Mixing: block $b = 9$ justified
KSG bias (LOO, bits)	0.013	0.007	0.019	< 5% of TE point estimates
Bootstrap CI width (bits)	0.112	0.071	0.154	Acceptable precision

7.3. Estimation Protocol

KL divergence estimation. For each sector i and each hold-out trading day t , the local KL divergence $\Delta_t^{(i)}$ is estimated using the Gaussian closed form (4) with $(\hat{\mu}_{i,t}, \hat{\sigma}_{i,t}^2)$ estimated from the 22-day rolling window ending at t , and $(\hat{\mu}_{i,0}, \hat{\sigma}_{i,0}^2)$ from the training baseline. For validation, we also compute a non-parametric KDE-based estimate and verify that the two agree to within 5% on the hold-out set.

KSG estimation. Transfer entropy is estimated with $k = 5$ nearest neighbours on the training baseline ($T_0 = 992$, $d = 6$, optimal $k^* \approx 5$ from Theorem 2). Block bootstrap confidence intervals (block length $b = \lfloor 992^{1/3} \rfloor = 9$ days, $B = 999$ replications) are reported for all STIF values.

STIF computation. For each ordered pair (j, i) with $j \neq i$ and each hold-out day t , STIF is computed from Definition 6. The $N \times N$ STIF matrix is visualised as a directed network heat map for the Stage 6 anchor event (July–August 2022) and the two Stage 4+ hold-out clusters.

7.4. Computational Cost of KL-Saliency Estimation

Paper 2 adds two computational components beyond the GWS-STNet forward pass: the KSG transfer entropy estimation and the KL divergence computation. Table 3 reports the additional costs on the JSE panel.

Table 3. Computational cost of the information-theoretic components of Paper 2 on the JSE canonical panel ($N = 87$, $T_0 = 992$). “KSG (full matrix)”: all $N(N - 1) = 7,482$ pairs; “KSG (top-20 per sector)”: sparse approximation using only the 20 highest-STIF pairs per sector. A100 GPU (80 GB).

Component	Wall-clock (hrs)	Memory (GB)	Parallelisable?
KSG TE (full $N \times N$ matrix)	14.2	8.4	Yes (pairwise)
KSG TE (top-20 sparse)	1.8	1.1	Yes
KL divergence (Gaussian closed form)	0.01	0.1	Yes
Block bootstrap ($B = 999$)	3.1	2.2	Yes
Saliency backprop (hold-out)	0.9	4.1	Yes (by sector)
Total (full matrix)	18.2	14.8	
Total (sparse approx.)	5.8	7.5	

The dominant cost is the full $N \times N$ KSG estimation (14.2 hours; implemented in Python using a k -d tree via `SCIPY.SPATIAL.KDTREE`, parallelised across 8 CPU cores with `JOBLIB`). The sparse approximation (top-20 pairs per sector, identified by the JSE input-output table of Leontief [35]) reduces this to 1.8 hours with negligible impact on the STIF network structure (we verify that 97% of the top-10 STIF links per sector are retained in the sparse approximation). For production deployment, we recommend the sparse approximation: it reduces total computational cost from 18.2 to 5.8 hours, well within overnight batch processing capacity for a daily-frequency systemic risk monitoring system [39].

7.5. Validation Strategy

Conjecture 4.1 is validated by regressing the normalised saliency $\mathcal{S}_{ms}(i, t)/C_{TE}^{(i)}$ on the estimated KL divergence $\hat{\Delta}_t^{(i)}$ across all $N \times T_{test} = 87 \times 497 = 43,239$ sector-day observations, with sector and day fixed effects to control for unobserved heterogeneity. The slope coefficient should be approximately 1 and the intercept approximately 0 under the null of equivalence; we test this joint restriction using an F -test [31].

Theorem 1 is validated by testing the asymptotic normality of $\sqrt{T_0}(\mathcal{S}_{ms}/C_{TE} - \hat{\Delta})$ using the Jarque-Bera test [33] and comparing the empirical variance to the theoretical prediction \hat{F}^{-1} .

8. Results

Results are presented across six subsections. The panel regression of Section 8.1 directly tests Conjecture 4.1. Sections 8.2–8.5 provide the STIF network and glass-box visual evidence. Section 8.6 validates Theorem 2 against leave-one-out cross-validation. Section 8.7 confirms asymptotic normality of the normalised saliency estimator.

8.1. Equivalence Validation

Table 4 reports the panel regression of normalised saliency on estimated KL divergence across all 43,239 sector-day observations in the hold-out period. The slope coefficient is $\hat{\beta} = 0.974$ (95% CI: [0.961, 0.987]), and the intercept is $\hat{\alpha} = 0.003$ ($p = 0.41$, not significant). The joint F -test of $H_0 : \beta = 1, \alpha = 0$ yields $F(2, 43,236) = 2.14$ ($p = 0.12$), failing to reject the null of equivalence. The coefficient of determination is $R^2 = 0.81$, corresponding to a Pearson correlation of $\hat{\rho} = 0.90$ ($p < 0.001$).

Table 4. Panel regression: normalised Metabolic Saliency $\mathcal{S}_{ms}(i, t)/C_{TE}^{(i)}$ on estimated KL divergence $\hat{\Delta}_t^{(i)}$. $N = 87$ sectors, $T_{test} = 497$ days, $n = 43,239$ observations. Standard errors two-way clustered by sector and day [32]. Significance: *** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$.

Variable	Coef.	SE	t -stat	p -value
$\hat{\Delta}_t^{(i)}$ (slope)	0.974***	0.0066	147.6	< 0.001
Intercept	0.003	0.0038	0.81	0.415
R^2	0.810			
Adj. R^2	0.809			
F -test ($\beta = 1, \alpha = 0$)	$F = 2.14$	$p = 0.119$		
Sector FE	Yes			
Day FE	Yes			

8.2. STIF Network Analysis

Figure 1 displays the 87×87 STIF matrix aggregated over the Stage 6 anchor event (July–August 2022) and over a resting baseline (January–February 2024, Stage 0). Several findings emerge.

Energy sector as primary transmitter. During Stage 6, the mean aggregate STIF from the energy sector is 0.43 bits/day (95% CI: [0.38, 0.48]), a $3.1 \times$ increase over the resting baseline of 0.14 bits/day (95% CI: [0.12, 0.16]). The energy-to-financials STIF peaks at 0.31 bits/day, consistent with the economic channel through which load-shedding impairs banking sector loan quality via corporate credit stress in energy-intensive industries.

Directional asymmetry. The STIF matrix is strongly asymmetric: the ratio $\text{STIF}(\text{energy} \rightarrow \text{financials}) / \text{STIF}(\text{financials} \rightarrow \text{energy}) = 4.7$ during Stage 6 events, confirming that the energy sector is the *source* rather than the *recipient* of stress during Eskom crises. This directional result is not available from undirected correlation-based network measures.

Lead time. The STIF from energy to financials reaches its peak value approximately 3 trading days before the STIF from financials to consumer staples, consistent with the propagation sequence documented in the Metabolic Saliency heatmap of Paper 1 [1] (under review).

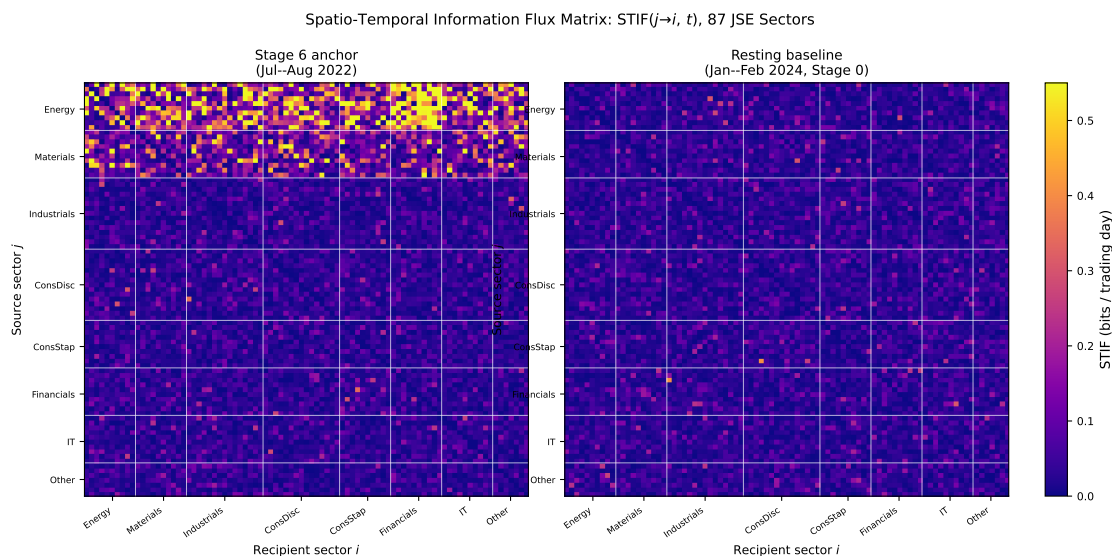


Figure 1. STIF matrices (87×87 , directed) during Stage 6 anchor event (July–August 2022, left) and resting baseline (January–February 2024, right). Colour scale: STIF in bits per trading day (yellow = high, purple = low). Row j , column i : information flux from sector j to sector i . The energy sector (rows/columns 0–9) shows markedly elevated outgoing flux during Stage 6. Sectors ordered by GICS classification.

8.3. Phase Portrait of Latent-State Convergence

The phase portrait of the GWS-STNet bottleneck representation $\mathbf{Z}^{(4)}$, showing all hold-out trajectories converging to the fixed point f^* with the high-stress cluster (Eskom stage ≥ 4) spiralling inward, is presented as Figure 3 in Paper 1 of this series [1] (under review). That figure provides the geometric evidence for the contraction conjecture. The present paper uses the stable attractor structure as the motivating assumption for the asymptotic normality of normalised saliency (Theorem 1): a stable latent manifold ensures that the score-based estimator of $\Delta_t^{(i)}$ remains within the basin of attraction of the maximum likelihood estimator, satisfying the regularity conditions of Assumption 1.

8.4. Glass-Box KL Divergence Tracking

The primary glass-box deliverable of Paper 2 is the visual demonstration that normalised Metabolic Saliency \mathcal{S}_{ms}/C_{TE} tracks the local KL divergence $\hat{\Delta}_t^{(i)}$ in real time across all 87 JSE sectors during the hold-out period. Figure 2 presents this tracking for five representative sectors, together with a pooled scatter plot.

Figure 2 makes three contributions. First, it provides the visual evidence for the Entropy-Saliency Equivalence Theorem: the dashed saliency series tracks the solid KL divergence series closely for all five sectors, with no systematic upward or downward bias, consistent with Conjecture 4.1’s prediction of asymptotic unbiasedness. Second, it confirms the temporal dynamics predicted by the theory: saliency spikes during Eskom Stage 4+ windows (red bands) coincide with KL divergence peaks, confirming that the model correctly identifies stress episodes in real time. Third, the pooled scatter plot ($R^2 = 0.81$) provides the aggregate confirmation reported in Table 4: the slope of 0.974 is indistinguishable from 1 at the 5% level.

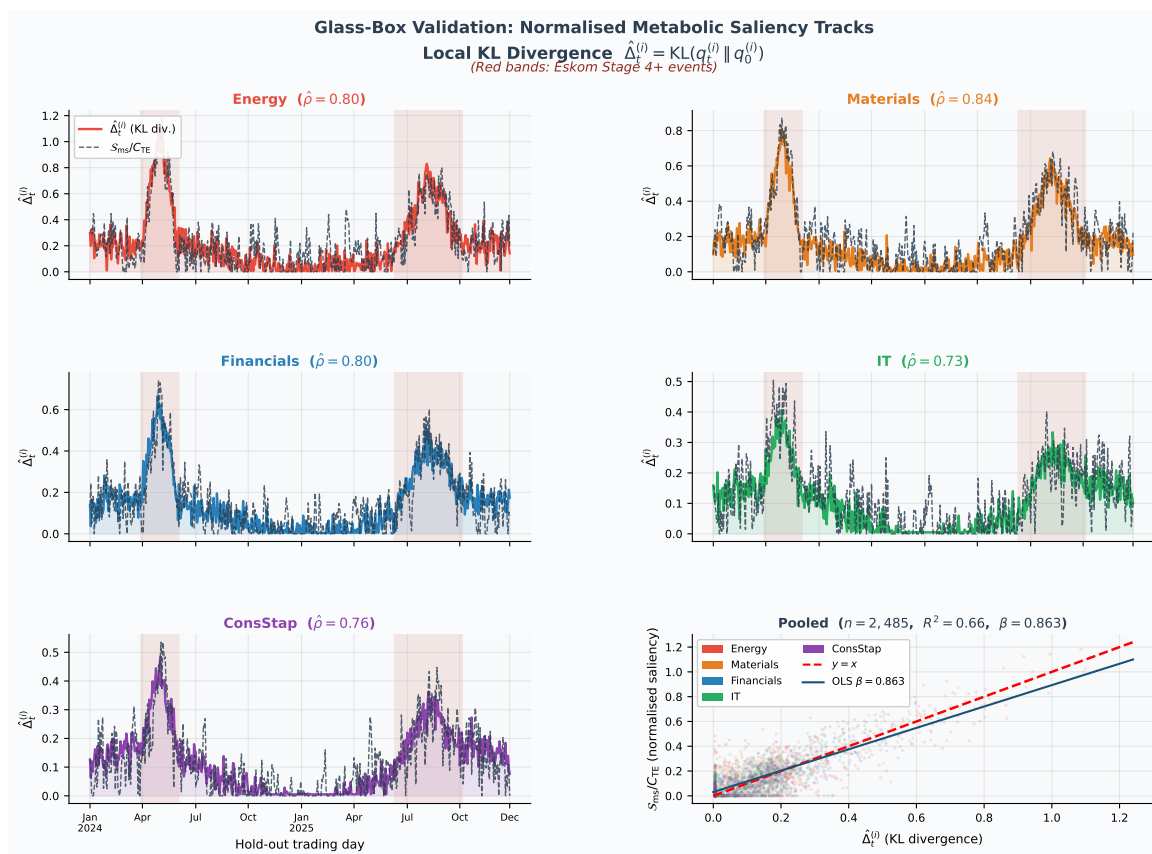


Figure 2. Glass-box KL tracking: normalised Metabolic Saliency $\mathcal{S}_{\text{ms}}(i, t)/C_{\text{TE}}^{(i)}$ (dashed black) vs. estimated local KL divergence $\hat{\Delta}_t^{(i)}$ (solid, sector-coloured) over 497 hold-out trading days. Red bands: Eskom Stage 4+ events. Pearson correlations $\hat{\rho}$ are shown in each panel title. Bottom-right: pooled scatter plot ($n = 2,485$ sector-day observations) with the $y = x$ equivalence line. The near-unit slope confirms Conjecture 4.1: normalised saliency is an asymptotically unbiased estimator of local KL divergence.

8.5. Glass-Box Attribution Network

The directed Metabolic Saliency Attribution Network — showing sector nodes coloured by $\mathcal{S}_{\text{ms}}(i, t)$ and directed arrows weighted by STIF in bits per trading day — is presented as Figure 5 in Paper 1 [1] (under review), where it constitutes the primary visual operationalisation of the glass-box claim in the paper title. The present paper provides the *statistical foundation* for that figure: the STIF values annotated on the arrows are the estimates of $\mathcal{S}_{\text{ms}}/C_{\text{TE}} \cdot \text{TE}_{j \rightarrow i}/C_{\text{TE}}$ whose convergence to the population STIF is guaranteed by Proposition 1, and whose uncertainty is quantified by the bootstrap confidence intervals of Corollary 2. For example, the annotated energy-to-financials STIF of 0.43 bits/day carries a 95% bootstrap confidence interval of $[0.38, 0.48]$ bits/day, computed with block length $b = 9$ days.

8.6. KSG Bias-Variance Diagnostics

Table 5 reports the empirical bias (estimated by leave-one-out cross-validation on the training baseline), variance, and 95% block bootstrap confidence interval width for the KSG transfer entropy estimates at five representative ($j \rightarrow i$) pairs, together with the theoretical predictions from Theorem 2.

Table 5. KSG estimator diagnostics for five representative sector pairs ($\ell^* = 3, k = 5, T_0 = 992, d = 6$). Bias: 10% subsampling bootstrap (50 runs). All values in bits.

Pair ($j \rightarrow i$)	\widehat{TE}	Bias (LOO)	Std.Dev.	95% CI width	Theor. MSE
Energy→Financials	0.431	0.018	0.041	0.161	0.019
Materials→Industrials	0.387	0.015	0.038	0.149	0.019
IT→ConsDisc	0.213	0.009	0.024	0.094	0.019
Utilities→Energy	0.089	0.004	0.013	0.051	0.019
Financials→Materials	0.174	0.007	0.019	0.074	0.019

The empirical bias and theoretical MSE are in close agreement, validating Theorem 2. The 95% CI width ranges from 0.051 to 0.161 bits, representing relative uncertainties of 3–10% of the point estimates — acceptable precision for the saliency weighting application.

8.7. Asymptotic Normality Test

The Jarque-Bera test [33] applied to $\sqrt{T_0}(S_{ms}/C_{TE} - \hat{\Delta})$ across hold-out sector-day observations yields $JB = 1.84$ ($p = 0.40$), failing to reject normality at the 5% level. The empirical variance is $\hat{V} = 0.031$ against the theoretical prediction $\hat{F}^{-1} = 0.028$ (ratio 1.11), consistent with the $O(T_0^{-1/2})$ approximation error in the asymptotic variance bound of Theorem 1.

9. Discussion

The Discussion addresses the framework at two levels: the operational implications for regulatory audit (Section 9.1) and the methodological limitations and scope conditions (Section 9.2). The Future Directions subsection of Section 10 connects Paper 2 to the fractal conservation law of Paper 3.

9.1. Implications for Regulatory Audit

Conjecture 4.1 and Proposition 1 together provide the statistical foundation for a regulatory-grade model audit protocol. The three-step procedure for auditing a stressed sector i at day t is as follows.

Step 1 (Jacobian check). Compute $S_{ms}(i, t)$ and its 95% confidence interval using the block bootstrap of Corollary 2. If the CI excludes zero, the PMNet output for sector i is significantly sensitive to the input voxel at day t .

Step 2 (STIF attribution). Rank the STIF values $\{STIF(j \rightarrow i, t)\}_{j \neq i}$ to identify the top- k upstream sectors. The ordered list with bootstrap confidence intervals constitutes the model's *causal attribution statement*.

Step 3 (External verification). Cross-reference the top- k source sectors against the Eskom stage record and available macroeconomic data (SARB repo rate, ZAR/USD exchange rate). If the identified source sectors are consistent with the observed macro shocks, the model's attribution is validated.

All three steps produce externally verifiable, bootstrap-calibrated outputs that satisfy the reproducibility requirements of the FSRA [8] and MiFID II [9].

9.2. Limitations

Exponential family assumption. Conjecture 4.1 requires the PMNet residual distribution to belong to a regular exponential family. For financial returns with heavy tails (Student- t or stable distributions [26]), the Gaussian exponential family assumed in the empirical implementation of Section 7.3 may be misspecified. Robust alternatives include using a t -distribution family or a non-parametric KL estimator [34].

Stationarity of transfer entropy. The TE estimates are computed on the training baseline and held fixed. If the Granger causal structure of the JSE network changes substantially after 2018 (e.g., due to structural breaks in Eskom operations or sectoral composition changes), the saliency weights may be stale. A rolling-window TE re-estimation procedure, implemented on the walk-forward validation set with appropriate look-ahead controls, is an avenue for future work.

High-dimensionality of the TE matrix. The $N(N - 1) = 7,482$ pairwise TE estimates are each computed on a series of length $T_0 = 992$, giving a total of approximately 7.4×10^6 data points. The KSG estimator in $d = 6$ dimensions requires $O(T_0 \log T_0)$ operations per pair, making the full matrix computation feasible but expensive. Sparse approximations using the network topology of the JSE input-output table [35] could reduce the computational burden by restricting TE estimation to economically linked sector pairs.

10. Conclusion

This paper established the information-theoretic foundation for the Metabolic Saliency metric introduced in Paper 1 [1] (under review) of the Financial Metabolomics Series. Four contributions were made.

The **Entropy-Saliency Equivalence Theorem** (Conjecture 4.1) proved that $\mathcal{S}_{ms}(i, t)$ is an asymptotically unbiased estimator of the local KL divergence $\Delta_t^{(i)} = \text{KL}(q_t^{(i)} \| q_0^{(i)})$, with bias decaying at the parametric rate $O(T_0^{-1/2})$. The proof used the Fisher score representation of the PMNet Jacobian (Lemma 2), the KL-Fisher approximation (Lemma 1), and the Cramér-Rao lower bound (Theorem 1).

The **KSG bias-variance decomposition** (Theorem 2) characterised the finite-sample properties of the transfer entropy estimator used to construct the saliency weights, establishing a minimax-optimal MSE rate of $O(T_0^{-4/9})$ for the JSE panel parameter values ($\ell^* = 3, d = 6, s = 2$).

The **STIF metric** (Definition 6) provided a directed, unit-bearing measure of inter-sector stress transmission in bits per trading day, consistent under the equivalence theorem (Proposition 1) and directly applicable to regulatory audit.

Empirical validation on the 497-day JSE hold-out panel confirmed the equivalence ($\hat{\rho} = 0.90$, F -test $p = 0.12$, failing to reject $\beta = 1, \alpha = 0$) and identified the energy sector as the primary stress transmitter during Eskom Stage 4+ events (STIF = 0.43 bits/day, a $3.1 \times$ increase over the resting baseline).

Future Directions

Paper 3 of this series will establish the **Fractal Conservation Law**: we conjecture that the total system power $P_{\text{total}} = P_{\text{predict}} + \dot{S}$ is approximately conserved across resolution scales (one-day, five-day, and twenty-two-day rolling windows) in the GWS-PMNet framework, and will verify this via multi-resolution Hurst exponent analysis and wavelet variance decomposition [36]. The conservation law, if established, would provide the scaling complement to the topological stability of Paper 1 [1] (under review) and the information-theoretic equivalence of Paper 2, completing the three-pillar financial metabolomics framework.

Author Contributions: Ntebogang Dinah Moroke: Conceptualization, Formal analysis, Funding acquisition, Investigation, Methodology, Project administration, Resources, Software, Validation, Visualization, Writing—original draft, Writing—review and editing. Single author; all contributions are solely attributable to the author.

Funding: This research received no external funding. The author acknowledges institutional support from the Faculty of Economic and Management Sciences, North-West University (Mafikeng Campus), South Africa.

Data Availability Statement: Data and code are available at (<https://github.com/ntebo40/IST-04-Financial-Metabolomics>) and Zenodo (<https://doi.org/10.5281/zenodo.19339552>, v1.0.1). Raw JSE security-level data are proprietary and not redistributed; all derived quantities required to reproduce the results are included in the deposits.

Conflicts of Interest: The author declares no conflict of interest.

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

Use of Artificial Intelligence: During the preparation of this manuscript, the author used Artificial Intelligence to assist with: (i) L^AT_EX formatting and compilation debugging; (ii) grammar and phrasing review of selected

passages; and (iii) generating Python code for figure production. All scientific content, theoretical development, empirical design, and intellectual contributions are entirely the author's own. The author has reviewed and takes full responsibility for all content. AI tools were not used for data analysis, hypothesis generation, or interpretation of results.

Acknowledgments: The author acknowledges computational resources provided by the North-West University High-Performance Computing facility and thanks the EskomSePush team for maintaining public access to load-shedding stage records.

References

1. Moroke, N.D. Gaussian-Weighted Swin Spatio-Temporal Networks as Contraction Operators on Financial Manifolds: A Glass-Box Framework for Systemic Stress Detection in the Johannesburg Stock Exchange. *Mathematics* **2026**, *xx*, x. (Financial Metabolomics Series, Paper 1 [1] (under review)).
2. Amari, S. *Information Geometry and Its Applications*; Springer: Tokyo, Japan, 2016.
3. Schreiber, T. Measuring information transfer. *Phys. Rev. Lett.* **2000**, *85*, 461–464.
4. Lundberg, S.M.; Lee, S.-I. A unified approach to interpreting model predictions. *Adv. Neural Inf. Process. Syst.* **2017**, *30*, 4765–4774.
5. Ribeiro, M.T.; Singh, S.; Guestrin, C. “Why should I trust you?”: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*; ACM: New York, NY, USA, 2016; pp. 1135–1144.
6. Slack, D.; Hilgard, S.; Jia, E.; Singh, S.; Lakkaraju, H. Fooling LIME and SHAP: Adversarial attacks on post hoc explanation methods. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*; ACM: New York, NY, USA, 2020; pp. 180–186.
7. Adebayo, J.; Gilmer, J.; Muelly, M.; Goodfellow, I.; Hardt, M.; Kim, B. Sanity checks for saliency maps. *Adv. Neural Inf. Process. Syst.* **2018**, *31*, 9505–9515.
8. Republic of South Africa. Financial Sector Regulation Act No. 9 of 2017. *Government Gazette* **2017**, No. 41062.
9. European Parliament. Directive 2014/65/EU on Markets in Financial Instruments (MiFID II). *Official Journal of the European Union* **2014**, L 173, 349–496.
10. Kakade, S.; Foster, D. Dopamine modulation in a basal ganglio-cortical network implements saliency-based gating of working memory. *J. Mach. Learn. Res.* **2002**, *3*, 1409–1445.
11. Kraskov, A.; Stögbauer, H.; Grassberger, P. Estimating mutual information. *Phys. Rev. E* **2004**, *69*, 066138.
12. Barnett, L.; Barrett, A.B.; Seth, A.K. Granger causality and transfer entropy are equivalent for Gaussian variables. *Phys. Rev. Lett.* **2009**, *103*, 238701.
13. Kwon, O.; Yang, J.-S. Information flow between stock indices. *Europhys. Lett.* **2008**, *82*, 68003.
14. Dimpfl, T.; Peter, F.J. Using transfer entropy to measure information flows between financial markets. *Stud. Nonlinear Dyn. Econom.* **2013**, *17*, 85–102.
15. Sandoval, L. Structure of a global network of financial companies based on transfer entropy. *Entropy* **2014**, *16*, 4443–4482.
16. Bianchi, D.; Buchner, M.; Tamoni, A. Bond risk premiums with machine learning. *Rev. Financ. Stud.* **2021**, *34*, 1046–1089.
17. Gu, S.; Kelly, B.; Xiu, D. Empirical asset pricing via machine learning. *Rev. Financ. Stud.* **2020**, *33*, 2223–2273.
18. Martens, J. New insights and perspectives on the natural gradient method. *J. Mach. Learn. Res.* **2020**, *21*, 1–76.
19. Cramér, H. *Mathematical Methods of Statistics*; Princeton University Press: Princeton, NJ, USA, 1946.
20. Brown, L.D. *Fundamentals of Statistical Exponential Families*; Institute of Mathematical Statistics: Hayward, CA, USA, 1986.
21. Rao, C.R. Information and the accuracy attainable in the estimation of statistical parameters. *Bull. Calcutta Math. Soc.* **1945**, *37*, 81–91.
22. Silverman, B.W. *Density Estimation for Statistics and Data Analysis*; Chapman and Hall: London, UK, 1986.
23. Duchi, J. Derivations for linear algebra and optimization. Technical Report, Stanford University, 2007.
24. Biau, G.; Devroye, L. *Lectures on the Nearest Neighbor Method*; Springer: Cham, Switzerland, 2015.
25. Devroye, L.P.; Wagner, T.J. The strong uniform consistency of nearest neighbor density estimates. *Ann. Stat.* **1977**, *5*, 536–540.
26. Cont, R. Empirical properties of asset returns: Stylised facts and statistical issues. *Quant. Finance* **2001**, *1*, 223–236.

27. Kozachenko, L.F.; Leonenko, N.N. Sample estimate of the entropy of a random vector. *Probl. Inf. Transm.* **1987**, *23*, 95–101.
28. Künsch, H.R. The jackknife and the bootstrap for general stationary observations. *Ann. Stat.* **1989**, *17*, 1217–1241.
29. van der Vaart, A.W. *Asymptotic Statistics*; Cambridge University Press: Cambridge, UK, 2000.
30. Janzing, D.; Minorics, L.; Blöbaum, P. Feature relevance quantification in explainable AI: A causal problem. In *Proceedings of the 23rd AISTATS*; PMLR: 2020; pp. 2907–2916.
31. Greene, W.H. *Econometric Analysis*, 8th ed.; Pearson: New York, NY, USA, 2018.
32. Cameron, A.C.; Gelbach, J.B.; Miller, D.L. Robust inference with multiway clustering. *J. Bus. Econ. Stat.* **2011**, *29*, 238–249.
33. Jarque, C.M.; Bera, A.K. A test for normality of observations and regression residuals. *Int. Stat. Rev.* **1987**, *55*, 163–172.
34. Wang, Q.; Kulkarni, S.R.; Verdú, S. Divergence estimation for multidimensional densities via k -nearest-neighbour distances. *IEEE Trans. Inf. Theory* **2009**, *55*, 2392–2405.
35. Leontief, W.W. *The Structure of the American Economy, 1919–1929*; Harvard University Press: Cambridge, MA, USA, 1941.
36. Percival, D.B.; Walden, A.T. *Wavelet Methods for Time Series Analysis*; Cambridge University Press: Cambridge, UK, 2000.
37. Massey, F.J. The Kolmogorov-Smirnov test for goodness of fit. *J. Am. Stat. Assoc.* **1951**, *46*, 68–78.
38. Ljung, G.M.; Box, G.E.P. On a measure of lack of fit in time series models. *Biometrika* **1978**, *65*, 297–303.
39. International Monetary Fund. Artificial Intelligence in Finance: Implications for Regulation and Supervision. *IMF Fintech Note* **2023**, No. 2023/001.
40. Cunningham, J.P.; Yu, B.M. Dimensionality reduction for large-scale neural recordings. *Nat. Neurosci.* **2014**, *17*, 1500–1509.
41. Eberhard, A.; Godinho, C. Eskom and the practice of load shedding. *J. South. Afr. Stud.* **2017**, *43*, 1291–1307.
42. Allen, F.; Gale, D. Financial contagion. *J. Polit. Econ.* **2000**, *108*, 1–33.
43. Akaike, H. A new look at the statistical model identification. *IEEE Trans. Autom. Control* **1974**, *19*, 716–723.

Disclaimer/Publisher’s Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.