

This is a preprint submitted to a journal. The final published version may differ.

Hybrid Context-Fusion Attention (CFA) U-Net and Clustering for Robust Seismic Horizon Interpretation

José Luis Lima de Jesus Silva^{‡}, João Pedro Gomes[†], Paulo Roberto de Melo Barros Junior[§], Vitor Hugo Serravalle Reis Rodrigues[¶], and Alexsandro Guerra Cerqueira^{*†}*

^{*}Oxaala Tecnologias

Rua Dinah Silveira de Queirós, Nº 06, Quinta do Candeal, Horto Florestal, Salvador, Bahia, CEP 40.296.160

[†]Geophysics Institute, Federal University of Bahia (UFBA)

R. Barão de Jeremoabo, Ondina, Salvador, Bahia, CEP 40170-290

[‡]Grupo de Estudo e Aplicação de Inteligência Artificial em Geofísica (GEAIG), Geophysics Institute, Federal University of Bahia (UFBA)

[§]Petrobras – Petróleo Brasileiro S.A.

Avenida República do Chile, 65, Centro, Rio de Janeiro, RJ, CEP 20031-912

[¶]Geological Survey of Brazil – Superintendência de Salvador

Av. Ulysses Guimarães, 2862, Sussuarana, Salvador, Bahia, CEP 41213-000

J. Luis Silva, J. Gomes, P. Barros Junior, V. Rodrigues, and A. Cerqueira

ABSTRACT

Interpreting seismic horizons is a critical task for characterizing subsurface structures in hydrocarbon exploration. Recent advances in deep learning, particularly U-Net-based architectures, have significantly improved automated horizon tracking. However, challenges remain in accurately segmenting complex geological features and interpolating horizons from sparse annotations. To address these issues, a hybrid framework is presented that integrates advanced U-Net variants with spatial clustering to enhance horizon continuity and geometric fidelity. The core contribution is the Context Fusion Attention (CFA) U-Net, a novel architecture that fuses spatial and Sobel-derived geometric features within attention gates to improve both precision and surface completeness. The performance of five architectures, the U-Net (Standard and compressed), U-Net++, Attention U-Net, and CFA U-Net, was systematically evaluated across various data sparsity regimes (10-, 20-, and 40-line spacing). This approach outperformed existing baselines, achieving state-of-the-art results on the Mexilhão field (Santos Basin, Brazil) dataset with a validation IoU of 0.881 and MAE of 2.49 ms, and excellent surface coverage of 97.6% on the F3 Block of the North Sea dataset under sparse conditions. The framework further refines merged horizon predictions (in-line and cross-line) using Density-Based Spatial Clustering of Applications with Noise (DBSCAN) to produce geologically plausible surfaces. The results demonstrate the advantages of hybrid methodologies and attention-based architectures enhanced with geometric context, providing a robust and generalizable solution for seismic interpretation in structurally complex and data-scarce environments.

INTRODUCTION

Seismic interpretation is essential for characterizing subsurface geological formations and plays a key role in the oil and gas industry (Mattos et al., 2021). A critical aspect of this process is identifying key horizons within amplitude volumes. These horizons represent reliable, continuous reflection surfaces characterized by stable wavelet signatures in seismic surveys. High-fidelity mapping improves geological interpretation, allowing analysis of amplitude variations that can reveal important geological features. Therefore, accurate identification of these horizons is also essential to unraveling the temporal dynamics and formative processes that shape geological structures.

As seismic interpretation advances, the need for comprehensive three-dimensional (3D) data sets has become more apparent, especially in regions with complex geological structures. Dorn (1998) emphasized the limitations of traditional interpretations based on two-dimensional (2D) profiles. In large-scale surveys comprising numerous inlines and crosslines, manual horizon interpretation becomes labor-intensive and time-consuming. To mitigate this, advances in auto-tracking technology have emerged as viable alternatives (Luo et al., 2023). These tools are capable of tracing reflections iteratively using seed points and structural cues. However, their performance can be limited in complex settings, such as areas with faults or chaotic reflectivity patterns (Wu et al., 2019).

To overcome these limitations, a wide range of strategies have been proposed to enhance the automation of horizon tracking. Marfurt et al. (1999) introduced waveform similarity, which serves as a proxy for geological continuity, while Stark (2003) proposed phase unwrapping to derive the Relative Geological Time (RGT), later extended with fault constraints by Wu and Zhong (2012). Other methods utilize slope-based metrics and Dynamic Time Warping (DTW) to correlate seismic traces (Hale, 2013; Wu and Fomel, 2018).

Initial attempts to integrate machine learning into horizon tracking explored multi-layer perceptrons (Harrigan et al., 1992; Kusuma and Fish, 1993). However, the transformative impact came with modern deep learning techniques. Powered by advances in GPU computing and large-scale data sets, Convolutional Neural Networks (CNNs) have yielded significant improvements in waveform classification and pixel-wise segmentation (LeCun et al., 2015; Wu and Zhang, 2018; Peters et al., 2019; Wu et al., 2019; Yang and Sun, 2020; Calhes et al., 2021; Ravasi and Birnie, 2022). Among these models, U-Net architectures (Ronneberger et al., 2015), with their encoder-decoder symmetry and skip connections, have proven particularly effective for segmentation of geological features (Wu et al., 2019).

Recent developments have expanded this paradigm to multi-scale 3D CNNs (Tschanen et al., 2020) and advanced U-Net variants such as U-Net++ (Zhou et al., 2018) and Attention U-Net (AlSalmi and Elsheikh, 2024). Despite progress, several critical challenges remain (Yu and Ma, 2021), such as the robust interpolation of geological horizons from sparse annotations (Poulinakis et al., 2023), and the precise segmentation of horizons in faulted and discontinuous environments, where signal ambiguity

and structural complexity degrade performance.

Some of these challenges have been addressed by introducing and validating a comprehensive hybrid framework for seismic horizon interpretation, combining deep learning with spatial clustering. To enhance the performance in geologically complex regions, a cross-attention mechanism has been integrated based on the Attention U-Net (Oktay et al., 2018), enabling the model to focus on salient seismic features. Additionally, a model-agnostic Density-based spatial clustering of applications with noise (DBSCAN) has been implemented as a post-processing step to filter out spurious and spatially incoherent predictions (Dhua et al., 2015), thereby improving the geological plausibility of reconstructed surfaces.

The contributions are threefold. First, the Context-Fusion Attention (CFA) U-Net is proposed, a novel architecture that augments attention gates with spatial and Sobel-based heads to improve geometric precision and surface completeness. Second, a systematic comparative evaluation of five variants of the U-Net (Standard U-Net, compressed U-Net, U-Net++, Attention U-Net, and CFA U-Net) is conducted under varying annotation sparsity (10, 20, and 40-line spacing) in two geologically distinct data sets: the F3 Block in the North Sea Graben and the faulted Mexilhão Field in the Santos Basin (Brazil). Third, the hybrid workflow, which combines attention mechanisms with DBSCAN-based clustering, yields more robust, accurate, and generalizable horizon interpretations than existing baselines.

The results confirm that the proposed Context-Fusion Attention (CFA) U-Net achieved the highest validation IoU (0.881) and lowest MAE (2.49 ms) on the faulted Mexilhão data set, outperforming all baselines. Furthermore, the model achieved the highest surface coverage (97.6%) in sparse data scenarios, surpassing even the high-recall performance of U-Net++. These results highlight a measurable trade-off between precision and completeness, and demonstrate that attention-based architectures, when enhanced with context fusion mechanisms, offer a superior inductive bias for seismic interpretation under both sparse and geologically complex conditions.

DATASET

To evaluate the performance of the methods on geologically diverse data, two distinct 3D seismic volumes were used, as shown in Figure 1. The volume of seismic data F3 (survey from the Dutch sector of the North Sea) includes 651 inline and 951 crossline sections with a spatial sampling of 25 meters. The time dimension spans 1,848 ms, sampled at an interval of 4 ms. The FS8 seismic horizon is characterized by plane-parallel high-amplitude reflectors representing Cenozoic-era marine transgression deposits (Schroot et al., 2005). The reservoirs in this block exhibit clear direct hydrocarbon indicator (DHI) signatures (Troccoli et al., 2022; Barbosa et al., 2024). The second data set (BS-400 block) covers the Mexilhão field in the Santos Basin (offshore of Brazil) which is renowned for its significant hydrocarbon reserves and natural gas. The Mexilhão field seismic volume presents a greater geological com-

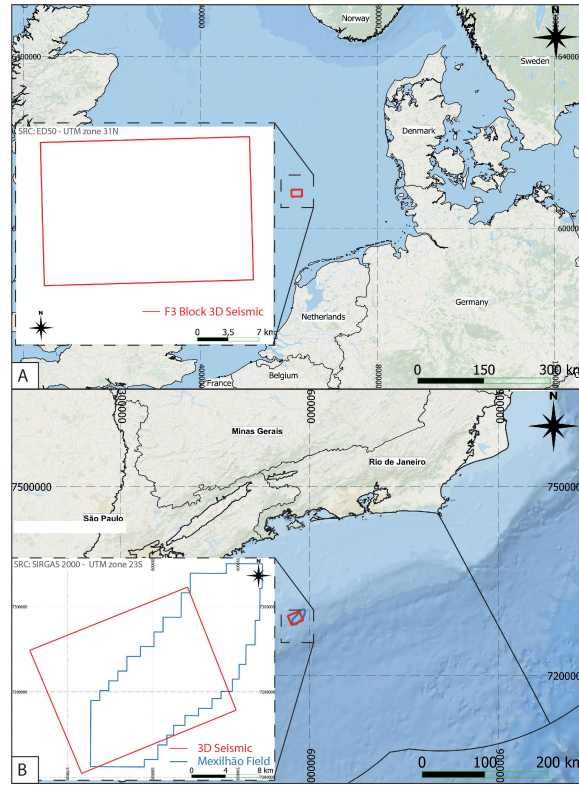


Figure 1: The two 3D seismic cubes utilized in the experiments are (a) the F3-Block in the Dutch North Sea and (b) the Mexilhão field in the Santos Basin, Brazil.

plexity, as illustrated in Figure 2. The data comprises 795 inlines and 624 crosslines, with a spatial resolution of 25 meters and a time range of 1600 ms, sampled at 4 ms. The geology is characterized by significant faulting resulting from halokinesis and extensional processes, which leads to reflectors with poor continuity in the rift section. Furthermore, the seismic data show good continuity in the reflectors of the drift portion. The largest reservoirs are turbiditic systems of the Ilha Bela Member, embedded within the shales of the Itajaí-Açu Formation. This complex structural and stratigraphic setting serves as an excellent test case for the robustness of the proposed segmentation framework.

To evaluate the models' ability to interpolate horizons from sparse interpretations, a systematic experiment was designed using both seismic cubes. For each of the six architectures, six distinct training scenarios were created by combining two orthogonal directions (inline and crossline) with three different sparsity levels, defined by the spacing between labeled lines (10, 20, or 40 lines). This resulted in a total of 36 unique trained models per seismic cube, allowing for a comprehensive assessment of each architecture's ability to generalize from sparse data. Table 1 specifies the number of labeled lines used for training in each configuration.

For each configuration, the data were partitioned using a systematic and non-random sampling strategy that directly simulates a real-world geophysical workflow.

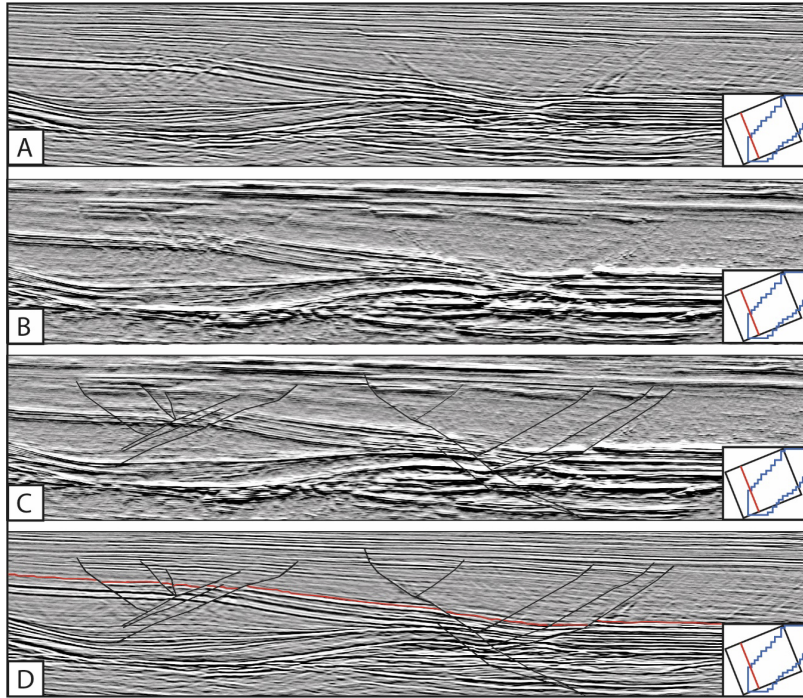


Figure 2: **Geological complexity of the Mexilhão field.** (a) A representative seismic amplitude slice. (b) The pseudo-relief attribute highlights structural features. (c) A manual fault interpretation overlaid on the pseudo-relief. (d) The seismic amplitude slice with both the fault interpretation and the target horizon highlighted.

data set	Line Spacing	# inlines	# crosslines
F3	10 x 10	96	47
	20 x 20	48	24
	40 x 40	24	12
Mexilhão	10 x 10	63	225
	20 x 20	32	113
	40 x 40	16	57

Table 1: Number of labeled inlines and crosslines used for training in each sparse data configuration for the two seismic volumes.

The sparsely sampled training set, for instance, consists of every 10th inline in the 3D volume. The remaining lines were reserved as the validation data set. From a machine learning perspective, this systematic partitioning serves as a form of spatial cross-validation Bergmeir and Benítez (2012), which is crucial to evaluate models on spatially correlated data such as seismic volumes (Roberts et al., 2017). This approach was deliberately chosen to avoid data leakage (Kaufman et al., 2012) that could occur with a standard random shuffle. In geophysics, adjacent seismic lines are usually highly correlated. Therefore, the model trained on line n and validated on the nearly

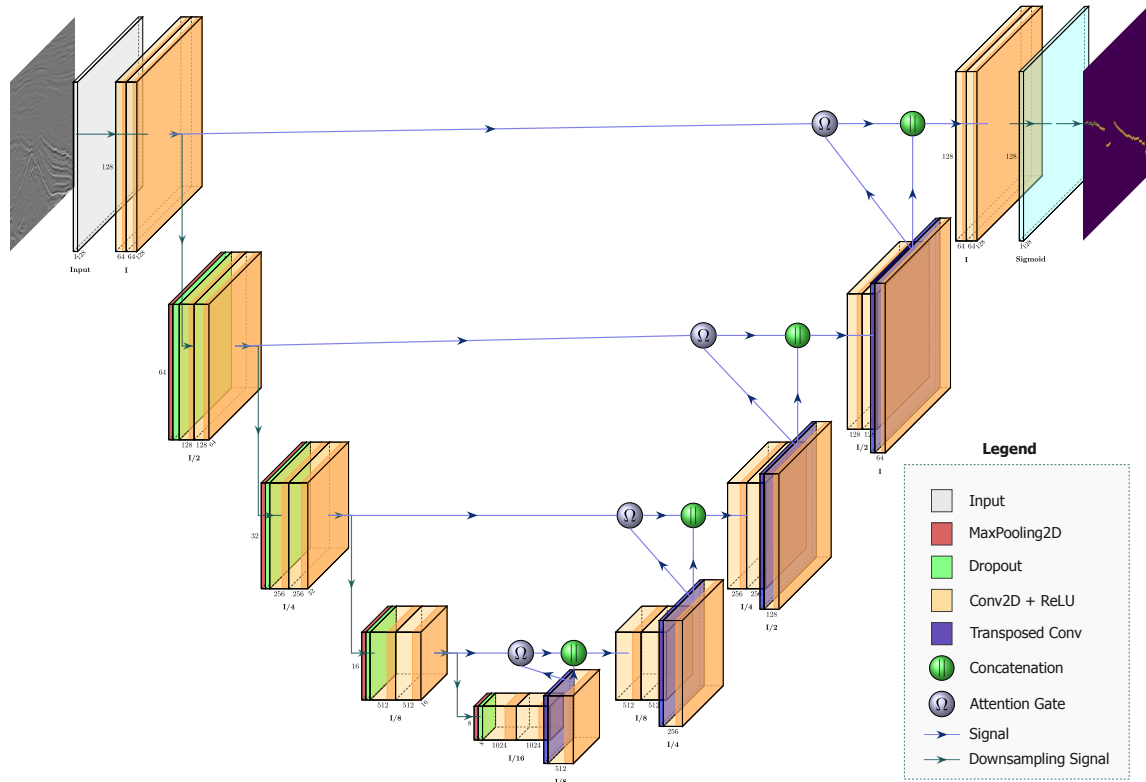


Figure 3: **Schematic of the Attention U-Net architecture.** The model processes a 2D seismic image patch (left) through a symmetric encoder-decoder network to produce a probability mask (right). The encoder (contracting path) uses convolutional (orange), dropout (green), and max-pooling (red) blocks to extract hierarchical features. The decoder (expansive path) uses upsampling (purple) and convolutional blocks to restore spatial resolution. The key components are the attention gates (Ω) on the skip connections, which adaptively re-weight features before they are fused via concatenation ($||$).

identical line $n + 1$ would yield artificially inflated performance metrics without truly testing the model’s ability to generalize. In contrast, this approach rigorously tests the model’s ability to interpolate geologically consistent horizons across large, unseen spatial gaps, which is the fundamental task required of an automated interpretation system.

METHODS

Architectural Framework

The segmentation of seismic horizons can be formulated as learning a mapping function $f : \mathbf{X} \rightarrow \mathcal{Y}$, and $\mathbf{X} \subset \mathbb{R}^{H \times W \times C}$ is the domain of 2D seismic image patches

with height H , width W , and C input channels, where $\mathcal{Y} \in \{0, 1\}^{H \times W}$ is the corresponding binary label space. To approximate this function, a family of deep convolutional neural networks derived from U-Net (Ronneberger et al., 2015) is investigated. This foundational model is distinguished by its symmetric encoder-decoder structure with long-range skip connections that fuse deep contextual features with shallow fine-grained details, making it exceptionally effective for precise localization. A systematic analysis of several U-Net variants with significant differences in parameterization is presented, particularly their efficacy when trained on sparsely labeled data sets common in geophysical applications.

To investigate the impact of model capacity, a compressed U-Net variant (Wu and Zhang, 2018) is included, which reduces the number of feature channels at each block. Furthermore, we explore two advanced architectures that modify the skip-connection mechanism. The U-Net++ architecture (Zhou et al., 2018) introduces nested and dense skip pathways to bridge the semantic gap between the encoder and decoder. The Attention U-Net (Oktay et al., 2018) integrates learnable attention gates into the skip connections to adaptively re-weight encoder features. A general architecture diagram is illustrated in Figure 3 using the Attention U-Net. The network follows a symmetric encoder-decoder structure to process 128×128 single-channel seismic patches, producing a probability mask of the same dimension. The encoder (contracting path) progressively downsamples the input through four blocks to capture high-level features, while the decoder (expansive path) symmetrically upsamples the feature maps to reconstruct the full-resolution output. The key innovation of the Attention U-Net, the attention gate (Ω), is integrated into each skip connection to filter features before they are concatenated ($||$) with the decoder path. The standard U-Net architecture is a simplification of this diagram where the attention gates are removed.

Based on these architectures, a key limitation in standard attention mechanisms has been identified, leading to the proposal of the Context Fusion Attention U-Net (*CFA U-Net*). The proposed model enhances the attention gate with a multi-head design that explicitly incorporates spatial and edge-aware inductive biases.

Mathematical Formulation of Key Architectures

As shown in the diagram Figure 3, the Attention U-Net modifies the standard U-Net architecture consisting of a symmetric encoder-decoder structure linked by skip connections that are modulated by attention gates.

The U-Net Encoder Pathway

The encoder path progressively extracts hierarchical features and reduces spatial dimensions. Given an input image $\mathbf{z}^0 \in \mathbb{R}^{H \times W \times C}$, the encoder comprises L levels. At each level $l \in \{1, \dots, L\}$, a convolutional block processes the input \mathbf{z}^{l-1} to produce a

feature map $\mathbf{x}_{\text{enc}}^l$, which is then passed to the corresponding decoder level via a skip connection, and then downsampled via max-pooling to produce the input for the next level (Krizhevsky et al., 2012), \mathbf{z}^l , is found using the following expressions:

$$\mathbf{x}_{\text{enc}}^l = \text{ConvBlock}_{\text{enc}}^l(\mathbf{z}^{l-1}) \quad (1)$$

$$\mathbf{z}^l = \text{MaxPool}(\mathbf{x}_{\text{enc}}^l) \quad (2)$$

where each $\text{ConvBlock}_{\text{enc}}^l$ consists of two sequential units, each comprising a 3×3 convolution, Batch Normalization (Ioffe and Szegedy, 2015), and a Rectified Linear Unit (ReLU) activation function (Nair and Hinton, 2010). The output of the final encoder level, \mathbf{z}^L , serves as the input to the bottleneck.

Attention Gated Skip Connections

To address the challenge of fusing semantically dissimilar features, Attention U-Nets (Oktay et al., 2018) employ attention gates (AGs) that adaptively re-weight the encoder feature maps to emphasize task-relevant features, $\mathbf{x}_{\text{enc}}^l$. The gating signal, \mathbf{g} , is the feature map from the next decoder level, \mathbf{d}^{l+1} , providing contextual information to guide the attention mechanism. The AG computes an attention coefficient map, $\alpha^l \in [0, 1]^{H_l \times W_l}$:

$$\alpha^l = \sigma_2 \left(\psi \left(\sigma_1 \left(\theta_x(\mathbf{x}_{\text{enc}}^l) + \theta_g(\text{Up}(\mathbf{g})) \right) \right) \right) \quad (3)$$

where $\theta_x(\cdot)$, $\theta_g(\cdot)$, and $\psi(\cdot)$ are linear projections implemented as 1×1 convolutions with their respective bias terms. The function σ_1 is a Rectified Linear Unit (ReLU) activation, and σ_2 is a sigmoid activation. The resulting attention-gated feature map, $\hat{\mathbf{x}}_{\text{enc}}^l$, is then computed via element-wise multiplication:

$$\hat{\mathbf{x}}_{\text{enc}}^l = \mathbf{x}_{\text{enc}}^l \odot \alpha^l \quad (4)$$

where $\hat{\mathbf{x}}_{\text{enc}}^l$ is the resulting **attention-gated feature map**, computed by performing an element-wise multiplication (\odot) between the original encoder feature map, $\mathbf{x}_{\text{enc}}^l$, and the attention coefficients, α^l . This operation effectively recalibrates the original features, emphasizing salient regions while suppressing irrelevant background information.

Architectural differences in Skip Pathways

In the standard U-Net (Ronneberger et al., 2015), the decoder receives features via a direct skip connection. At each level, the input to the convolutional block, $\text{ConvBlock}_{\text{dec}}^l$, is formed by concatenating the up-sampled feature map from the level below, $\text{Up}(\mathbf{d}^{l+1})$, directly with the original encoder feature map, $\mathbf{x}_{\text{enc}}^l$. This can be viewed as a simplification of the attention-based decoder, Equation 6, where the modulated feature map, $\hat{\mathbf{x}}_{\text{enc}}^l$, is replaced by the unmodified encoder features. This direct

fusion combines deep, semantic information from the decoder path with shallow, high-resolution spatial details from the encoder path at each level.

The U-Net++ architecture (Zhou et al., 2018), conversely, redesigns the skip pathways to be nested and densely connected, aiming to bridge the semantic gap between the encoder and decoder feature maps more effectively. This is achieved by creating a series of intermediate convolutional blocks along the skip connections. Let $\mathbf{x}^{l,j}$ denote the output of the node at down-sampling level l and convolutional layer j of the dense skip path, where $l = 0$ represents the top level and $j = 0$ represents the output of the encoder. The output of any node $\mathbf{x}^{l,j}$ for $j > 0$ is computed by aggregating the outputs of all previous nodes at the same level, along with the up-sampled output from the corresponding node in the level below:

$$\mathbf{x}^{l,j} = \text{ConvBlock}^{l,j} \left(\text{concat} \left[\left(\mathbf{x}^{l,k} \right)_{k=0}^{j-1}, \text{Up}(\mathbf{x}^{l+1,j-1}) \right] \right) \quad \text{for } j > 0 \quad (5)$$

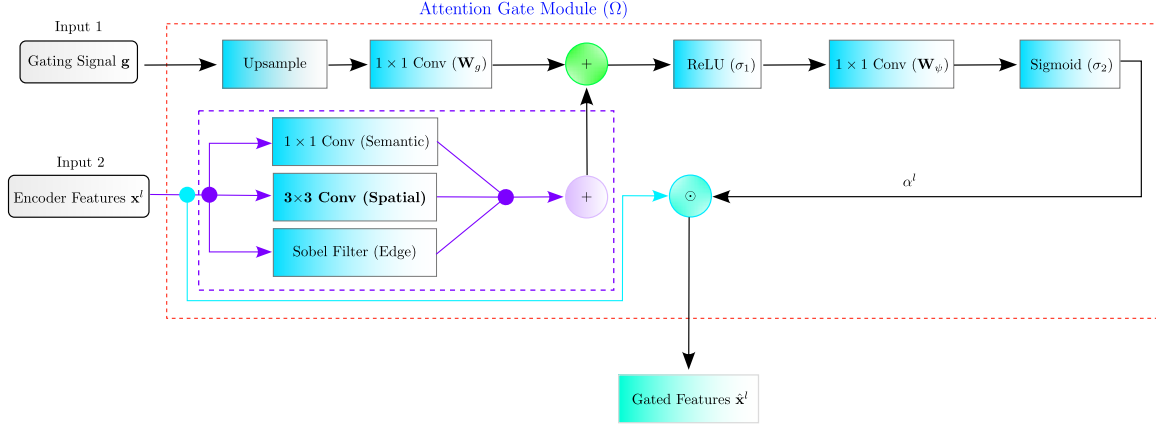
where $\text{ConvBlock}^{l,j}$ is the convolutional block at node (l, j) , consisting of the same operations (convolution, normalization, activation) as those in the encoder pathway, and $\text{concat}[\cdot]$ is the channel-wise concatenation operation. This formulation replaces the direct skip connection with a redesigned pathway where the decoder receives a densely aggregated set of feature maps. These maps are progressively enriched, facilitating a more gradual fusion of semantic and spatial information.

Decoder Path

The decoder path symmetrically restores the spatial resolution of the feature maps to that of the input image, while integrating the high-resolution, context-aware features from the attention-gated skip connections. The process begins with the feature map produced by the bottleneck, which we denote by \mathbf{d}^{L+1} . The decoder then proceeds through L levels, from $l = L$ down to 1. At each level l , the feature map from the previous (deeper) layer, \mathbf{d}^{l+1} , is first up-sampled by the operator $\text{Up}(\cdot)$. Upsampling is implemented as a learned 2×2 transposed convolution (Zeiler and Fergus, 2014), which allows the network to recover spatial information. The resulting tensor is then concatenated with the corresponding attention-gated feature map of the encoder, $\hat{\mathbf{x}}_{\text{enc}}^l$. This combined map is subsequently processed by a convolutional block of the decoder, $\text{ConvBlock}_{\text{dec}}^l$, to produce the output feature map \mathbf{d}^l :

$$\mathbf{d}^l = \text{ConvBlock}_{\text{dec}}^l \left(\text{concat} \left[\text{Up}(\mathbf{d}^{l+1}), \hat{\mathbf{x}}_{\text{enc}}^l \right] \right) \quad (6)$$

After the final decoder stage (at $l = 1$), a 1×1 convolution is applied. This final layer acts as a pixel-wise classifier, mapping the multi-channel feature representation from the last decoder block to a single logit for each pixel. A sigmoid activation function is then applied to this output to produce the final probability map (horizon target) $\hat{\mathbf{y}} \in [0, 1]^{H \times W \times 1}$.

Figure 4: **Diagram of Context-Fusion Attention gate mechanism**

Context-Fusion Attention U-Net (CFA U-Net)

Although the standard attention gate effectively reweighs features, its dependence on 1×1 convolutions makes it inherently spatially unaware. It determines the salience of features based on channel-wise information and cannot directly process local spatial patterns. This is a notable limitation for seismic interpretation, where key features are primarily defined by their local structure. To address this, the context-fusion attention gate, illustrated in Figure 4, was designed to be sensitive to semantic content, learned spatial patterns, and explicit edge characteristics simultaneously.

The core innovation is the replacement of the standard encoder feature projection, $\theta_x(\mathbf{x}_{\text{enc}}^l)$ in Equation 3, with an enriched, multi-head representation, $\mathbf{h}_{\text{fused}}$. This is created by fusing the outputs of three parallel heads, where each acts as a specialized feature extractor:

First, a **Semantic Head** preserves the standard 1×1 convolution to identify the channel-wise importance of statistical properties, such as high seismic amplitude.

$$\mathbf{h}_{\text{sem}} = \theta_{\text{sem}}(\mathbf{x}_{\text{enc}}^l) \quad (7)$$

Second, a **Spatial Head** employs a 3×3 convolution to introduce a crucial **spatial inductive bias**, an architectural assumption that neighboring pixels are structurally related. This bias allows the model to recognize geometric forms, such as the continuity of a reflector, even where its statistical signature is weak.

$$\mathbf{h}_{\text{spatial}} = \theta_{\text{spatial}}(\mathbf{x}_{\text{enc}}^l) \quad (8)$$

Finally, an **Edge-Aware Head** provides a strong, **task-relevant prior** by incorporating a fixed-weight Sobel filter. By injecting explicit gradient information, this head frees the model from having to learn fundamental edge detection, allowing it

to focus its capacity on learning the significance of those edges for the segmentation task.

$$\mathbf{h}_{\text{edge}} = \theta_{\text{edge}}(\text{Sobel}(\mathbf{x}_{\text{enc}}^l)) \quad (9)$$

The outputs of these three heads from Equations 7, 8, and 9 are then fused via element-wise addition to create the following representation:

$$\mathbf{h}_{\text{fused}} = \mathbf{h}_{\text{sem}} + \mathbf{h}_{\text{spatial}} + \mathbf{h}_{\text{edge}} \quad (10)$$

This fused tensor, $\theta_{x,\text{fusion}} = \mathbf{h}_{\text{fused}}$, directly replaces the $\theta_x(\mathbf{x}_{\text{enc}}^l)$ term in the attention mechanism. The attention coefficients for the CFA gate, α_{cfa}^l , are therefore computed by this direct substitution into Equation 3, such that:

$$\alpha_{\text{cfa}}^l = \sigma_2(\psi(\sigma_1(\mathbf{h}_{\text{fused}} + \theta_g(\text{Up}(\mathbf{g})))) \quad (11)$$

By endowing the attention gate with a multi-headed receptive field, the network learns not only *what* features are relevant but also *how they are structurally arranged*, significantly improving the segmentation of geologically complex regions. This design overcomes the spatial unawareness of standard attention gates. By creating three parallel processing streams, the fused knowledge enables the attention gate to make a more informed decision about not only which features are relevant but also how they are structurally arranged.

Optimization and Evaluation

To optimize the network parameters, a composite loss function \mathcal{L} that combines the Dice Loss ($\mathcal{L}_{\text{Dice}}$) and Binary Cross-Entropy (\mathcal{L}_{BCE}) was employed. The BCE loss for a single pixel is given by:

$$\mathcal{L}_{\text{BCE}}(y_i, \hat{y}_i) = -[y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] \quad (12)$$

where \hat{y}_i is the predicted probability for pixel i and y_i is the corresponding ground truth label (0 or 1). The Dice Loss is particularly effective for handling class imbalance in segmentation tasks and is defined as:

$$\mathcal{L}_{\text{Dice}} = 1 - \frac{2 \sum_{i=1}^N y_i \hat{y}_i + \epsilon}{\sum_{i=1}^N y_i^2 + \sum_{i=1}^N \hat{y}_i^2 + \epsilon} \quad (13)$$

where N is the total number of pixels, and ϵ is a small constant for numerical stability. The final loss function is a weighted sum, $\mathcal{L} = \alpha \mathcal{L}_{\text{BCE}} + \beta \mathcal{L}_{\text{Dice}}$, with parameters $\alpha = 0.5$ and $\beta = 0.5$. Model performance was evaluated using the Intersection over Union (IoU) metric, or Jaccard index, defined as:

$$\text{IoU} = \frac{|\mathbf{Y} \cap \hat{\mathbf{Y}}|}{|\mathbf{Y} \cup \hat{\mathbf{Y}}|} = \frac{\text{TP}}{\text{TP} + \text{FP} + \text{FN}} \quad (14)$$

Where \mathbf{Y} and $\hat{\mathbf{Y}}$ are the ground truth and predicted segmentation masks, and TP, FP, and FN are the true positives, false positives, and false negatives, respectively.

DBSCAN-based Filtering of Horizon Predictions

To refine the raw output from the neural networks and remove spurious, high-probability voxels that are geologically inconsistent with a continuous horizon, we employ the Density-Based Spatial Clustering of Applications with Noise (DBSCAN) algorithm as a postprocessing step (Ester et al., 1996; Deng, 2020). DBSCAN is exceptionally well-suited for this task as it can identify clusters of arbitrary shape and effectively isolate noise points in low-density regions without requiring a predefined number of clusters (Schubert et al., 2017).

The behavior of the DBSCAN algorithm is governed by two key parameters: the maximum neighborhood distance, ϵ , and the minimum number of points required to form a dense region, *MinPts*. Based on these parameters, each point in the data set is categorized according to its density. A point is classified as a **core point** if at least *MinPts* neighbors are found within its ϵ radius. A **border point**, in contrast, is not a core point itself but falls within the ϵ -neighborhood of one. Any point that satisfies neither of these conditions is subsequently labeled as a **noise point**. The algorithm forms final clusters by identifying connected components of core points and their associated border points, while all points designated as noise are discarded from the final result.

The input to the DBSCAN algorithm is a 3D point cloud generated from the probability volume, $\hat{\mathbf{y}}$, predicted by each architecture. To create the point cloud, a probability threshold of $\tau = 1 \times 10^{-5}$ was first applied to the volume. The spatial coordinates (i, j, k) of all suprathreshold voxels, where $\hat{y}_{ijk} > \tau$, were then extracted to form the point cloud for clustering.

The selection of the DBSCAN hyperparameters, ‘epsilon’ (ϵ) and ‘MinPts’, was determined empirically to suit the density of the predicted seismic data. The ‘MinPts’ value was set to 25 to ensure that only geologically significant point groupings were considered as clusters. The ‘epsilon’ parameter was calculated based on a vertical exaggeration factor ($z_{\text{factor}} = 3$), which results in an effective neighborhood distance of $\epsilon = 6.0$ for clustering. After applying DBSCAN with these parameters, the algorithm identifies all clusters in the data. The final step of the filtering process involves retaining only the single most significant point cluster, which is presumed to be the actual horizon, and discarding all smaller clusters and noise points.

Hybrid Segmentation Workflow

The methodology for seismic horizon segmentation follows a multistage process that integrates deep learning for initial prediction with density-based clustering for post-processing and refinement, as illustrated in Figures 5, 6, and 7. The DBSCAN filtering workflow is composed of three steps, as shown in Figure 5. First, **(a)** the raw 3D probability volume is displayed as a point cloud where the color represents the two-way travel time (TWT), ranging from shallow (yellow) to deep (purple). However,

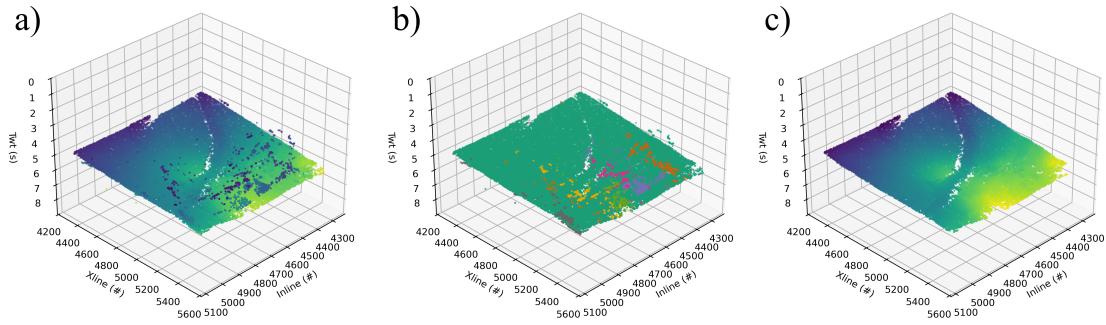


Figure 5: **The DBSCAN filtering workflow.** (a) The raw 3D probability volume is predicted by the neural network, with points colored by their two-way travel time (TWT). (b) DBSCAN groups the data into distinct clusters, identified by different categorical colors. (c) The final result after filtering for the largest point cluster, with the TWT-based color map reapplied to the cleaned horizon.

the prediction contains significant and spatially isolated noise. Therefore, **(b)** the points are clustered by distinct categorical colors, where the main horizon becomes a large cluster (green) and the noise is separated into smaller clusters of different colors. Finally, **(c)** only the most significant cluster is retained, producing a clean surface where the TWT color map is reapplied to show the final geologically coherent structure. This postprocessing is part of the final stage of the hybrid workflow, which is composed of three primary stages: experimental design for training on sparse data, model implementation and training, and the inference and filtering pipeline.

A systematic evaluation of the models was performed under varying data sparsity conditions. For each of the six architectures, six different training scenarios were created by combining two orthogonal interpretation directions (inline and crossline) with three label spacings (10th, 20th, or 40th labeled line). This experimental design resulted in a total of 36 unique trained models, allowing a comprehensive assessment of the ability of each architecture to interpolate between sparsely labeled seismic lines. A systematic, nonrandom split was employed in which using every 10th, 20th, or 40th line resulted in training sets that comprise approximately 10%, 5%, and 2.5% of the total available data, respectively. The remaining 90%, 95%, and 97.5% lines, respectively, were retained for validation and testing.

All models were implemented in Tensorflow/Keras and trained on an NVIDIA 4080 GPU using the Adam optimizer (Kingma and Ba, 2014). While a common procedure was followed, key hyperparameters such as learning rate and batch size were tuned for each architecture based on performance on a dedicated validation set, with the following learning rate (LR) and batch size (BS) pairs: **U-Net** (LR: 1×10^{-4} , BS: 1), **U-Net++** (LR: 5×10^{-3} , BS: 1), **Compressed U-Net** (LR: 5×10^{-4} , BS: 5), and **Attention U-Net** (LR: 5×10^{-4} , BS: 1). For the compressed U-Net, L2 regularization with a factor of 1×10^{-4} was applied to convolutional kernels to mitigate overfitting. Models were trained for up to 500 epochs using the DICE loss function,

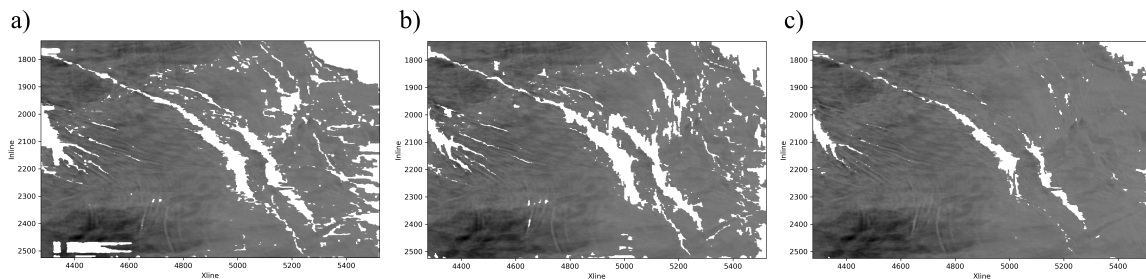


Figure 6: **Orthogonal prediction fusion.** This top-down view shows that predictions from models trained on (a) inline and (b) crossline directions are complementary. (c) represents the final surface with significantly improved completeness.

and an early stopping callback that monitored validation loss and ceased training if no improvement was observed for 30 consecutive epochs. The architectural framework was designed to be flexible for input tensor dimensions, allowing models to be adapted to different data set resolutions through a programmatic shape optimization utility.

After training, each model generated a 3D probability volume, $\hat{\mathbf{y}}$, which was refined using a two-stage postprocessing pipeline. The DBSCAN was applied to filter the raw probability volume and retain only the most significant and densely connected cluster of points. This step effectively removes spurious predictions that are spatially isolated from the main horizon structure, as illustrated in Figure 5. To maximize spatial coverage and mitigate directional biases, the final step fused the filtered predictions from models trained on orthogonal inline and crossline data, which is an effective technique in volumetric image analysis (Stollenga et al., 2015). The fusion is performed by taking the set union of the two filtered 3D point clouds, as illustrated in Figure 6:

$$\mathbf{P}_{\text{final}} = \mathbf{P}_{\text{inline}} \cup \mathbf{P}_{\text{xline}} \quad (15)$$

where \mathbf{P} represents a filtered point cloud for a given model and spacing. This process combines the strengths of both orthogonal predictions to produce a more complete and geologically coherent horizon surface.

Figure 7 shows the predicted horizon surfaces colored by the two-way travel time (TWT) (Sheriff and Geldart, 1995) (red is shallow, purple is deep). Using the U-Net++ architecture as an example, the figure highlights the layout that highlights the complementary nature of the orthogonal predictions. The rows compare the separate inline and crossline predictions against the final merged surface. Furthermore, the columns show a progressive filtering of the surface by depth, a technique that simplifies the identification of gaps and discontinuities.

The models' complementary performance is a direct result of interpreting three-dimensional, often anisotropic, geology using 2D seismic slices. For instance, the inline-trained model accurately captures the right side of the volume where the crossline model is weaker, while the crossline model resolves gaps in the upper surface that the inline model misses. By merging the two predictions via a set union

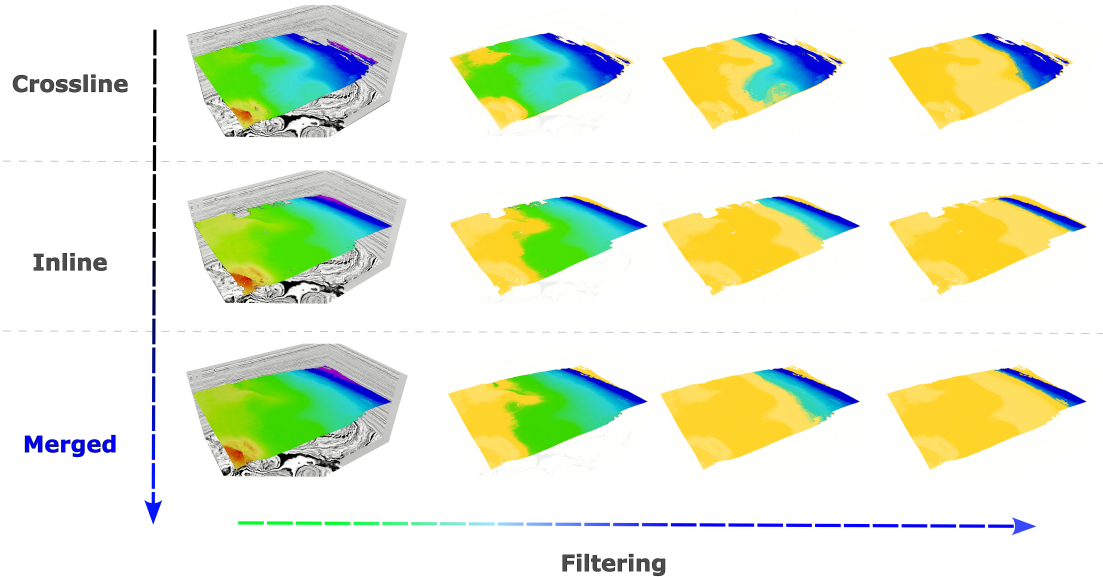


Figure 7: **Qualitative analysis of the orthogonal fusion workflow.** The top two rows show complementary predictions from models trained on crossline and inline directions, respectively. The bottom row shows the merged union, which produces a more complete horizon. The horizontal axis represents a progressive filtering of the surface by depth to visualize continuity.

given by Equation 15, the distinct strengths of each view-dependent interpretation are combined. Because a geological feature can have a clear signature in one view but appear as an ambiguous pattern in another, the merged horizon is more complete and geologically robust than an interpretation from a single orientation.

RESULTS

F3-Block Segmentation Results

The initial evaluation was performed on the F3 seismic volume, using the well-defined FS8 horizon as the segmentation target. Table 2 presents the mean training and validation metrics for each architecture across the three data sparsity configurations. All models achieved high accuracy scores (often > 0.99), however, this metric can be misleading in segmentation tasks with severe class imbalance, where the background (non-horizon) class constitutes the vast majority of pixels (Jadon, 2020). Therefore, the Intersection over Union (IoU) and Dice Loss provide more robust measures of performance. The results indicate that the standard U-Net architecture has strong and consistent performance, achieving both the highest validation IoU and the lowest

Architecture	Spacing	Mean Acc.		Mean IoU		Mean Dice Loss	
		Train	Valid	Train	Valid	Train	Valid
U-Net (Comp.)	10	0.9979	0.9942	0.9635	0.9039	0.3280	0.1727
	20	0.9971	0.9936	0.9496	0.8947	0.1343	0.1890
	40	0.9961	0.9901	0.9312	0.8461	0.2472	0.2364
U-Net	10	0.9992	0.9970	0.9843	0.9458	0.0278	0.0741
	20	0.9989	0.9970	0.9803	0.9464	0.0332	0.0740
	40	0.9981	0.9958	0.9653	0.9268	0.0614	0.0792
U-Net++	10	0.9995	0.9963	0.9901	0.9331	0.0218	0.0876
	20	0.9993	0.9967	0.9872	0.9403	0.0289	0.0816
	40	0.9990	0.9953	0.9820	0.9181	0.0603	0.0941
Attn. U-Net	10	0.9996	0.9954	0.9931	0.9177	0.0197	0.1006
	20	0.9993	0.9946	0.9865	0.9032	0.0288	0.1337
	40	0.9990	0.9916	0.9812	0.8510	0.0712	0.1990
CFA^S U-Net	10	0.9995	0.9951	0.9897	0.9127	0.0224	0.1174
	20	0.9992	0.9947	0.9859	0.9056	0.0279	0.1214
	40	0.9989	0.9917	0.9795	0.8571	0.1111	0.2360
CFA U-Net	10	0.9989	0.9966	0.9796	0.9381	0.0332	0.0843
	20	0.9980	0.9942	0.9631	0.8974	0.0521	0.1401
	40	0.9983	0.9952	0.9677	0.9165	0.0650	0.1016

Table 2: **Segmentation metrics on the F3 data set.** "Valid" refers to the validation set. Values represent the mean of models trained on inline and crossline data. IoU (higher is better) and Dice Loss (lower is better) are the most informative metrics for this unbalanced segmentation task.

validation Dice Loss across all data spacings, peaking with an IoU of 0.946 and a loss of 0.074 at the 10- and 20-line spacings. The proposed CFA **U-Net** also demonstrated highly competitive performance, ranking second in validation IoU on the 10-line spacing data set with a validation IoU of 0.938.

Following the pixel-level evaluation, the practical utility of each model is framed by the classic precision-recall trade-off (Davis and Goadrich, 2006), with results summarized in Table 3. Precision is defined as geometric accuracy, measured by the Mean Absolute Error (MAE) and Mean Squared Error (MSE). Lower values for these metrics indicate a more precise reconstruction, with predicted points being more faithful to the actual reflector depth. Recall is defined as spatial completeness, measured by the percentage of surface coverage area. Higher values indicate a more complete reconstruction of the entire surface. This framework allows for a precise analysis of how each architecture manages the inherent tension between these competing objectives.

The standard U-Net and U-Net++ architectures operate as high-recall models, consistently producing more complete surfaces with coverage areas in the range 97-98% (Table 3). However, this high recall can be associated with higher geometric

Architecture	Spacing	MAE	MSE	Area (%)
U-Net (Comp.)	10	4.21	22.69	96.34
	20	4.20	22.69	95.27
	40	4.35	25.86	91.74
U-Net	10	4.23	23.26	98.27
	20	4.25	24.00	98.20
	40	4.38	28.71	97.40
U-Net++	10	4.19	22.38	97.92
	20	4.22	23.06	98.11
	40	4.27	24.06	97.32
Attn. U-Net	10	4.14	21.39	93.26
	20	4.18	22.31	91.63
	40	4.21	23.04	86.77
CFA^S U-Net	10	4.77	45.97	96.83
	20	5.05	67.30	96.17
	40	5.18	54.27	93.16
CFA U-Net	10	4.44	33.99	98.19
	20	4.65	43.11	95.92
	40	4.83	47.83	97.61

Table 3: This table shows the geometric errors (MAE, MSE) and surface coverage area of the merged and filtered results for each model in the F3 data set. The proposed models are CFA^S U-Net (semantic and spatial heads) and the CFA U-Net (semantic, spatial, and Sobel heads).

errors (MAE of 4.38 for U-Net and 4.21 for Attn. U-Net at 40-line spacing). In contrast, Attention U-Net (Attn. U-Net) demonstrates a high precision, low recall strategy, consistently achieving the lowest MAE (ranging from 4.14 to 4.21) and MSE across all sparsity levels, indicating a high degree of geometric precision for the points it predicts. This precision comes at the cost of significantly lower surface coverage (recall), which drops to a low of 86.77% at 40-line spacing. This behavior is qualitatively confirmed in Figure 8, which shows larger gaps in the surfaces predicted by the Attention U-Net.

Furthermore, the standard U-Net demonstrates superior robustness to data sparsity. As training data decreases from 10- to 40-line spacing, its validation IoU remains remarkably stable, dropping by only 0.019 (from 0.946 to 0.927), while its surface coverage barely changes, declining by less than one percentage point (from 98.27% to 97.40%). In contrast, the more complex models show significant performance degradation. The Attention U-Net’s validation IoU drops more steeply by 0.067 (from 0.918 to 0.851). This represents more than three times the degradation of the standard U-Net, while its coverage falls sharply by 6.5 percentage points (93.26% to 86.77%). The CFA^S U-Net exhibits a similar drop in IoU, falling from 0.913 to 0.857. This sug-

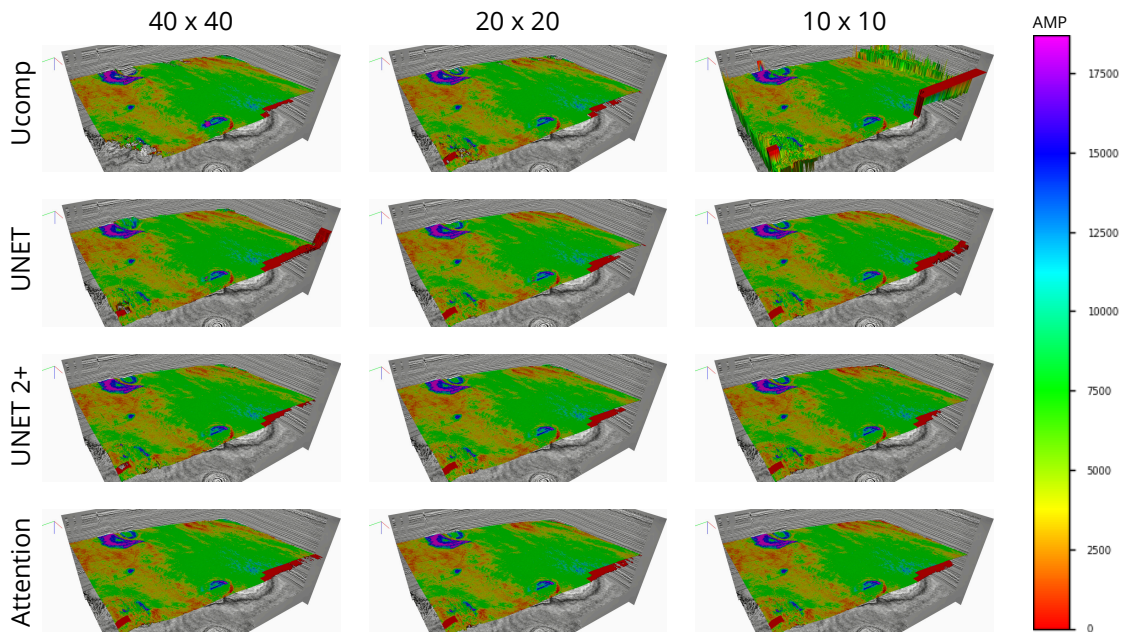


Figure 8: **Reconstructed 3D horizon surfaces after DBSCAN filtering for all architectures and data sparsity levels.** Each panel displays the final surface generated by a specific architecture (rows), trained on sparse grids with 40-line, 20-line, and 10-line spacing (columns). The figure visually demonstrates a key performance trade-off: while the standard U-Net and U-Net++ produce more complete surfaces, the Attention U-Net (Att-UNet) renders less complete horizons with larger gaps.

gests the simpler convolutional architecture of the standard U-Net possesses stronger generalization capabilities when data is limited.

The proposed Context-Fusion Attention model was designed to address the characteristic high-precision, low-recall behavior of the standard Attention U-Net. The *CFA^S U-Net* incorporates a learned 3×3 spatial head to encourage interpolation, while the full *CFA U-Net* adds a Sobel filter head to provide explicit edge guidance. As shown in Table 3, both models successfully improved surface coverage (recall), especially in sparse data conditions. At 40-line spacing, where the standard Attention U-Net coverage drops to 86.77%, the *CFA^S U-Net* increases to 93.16%, a relative improvement of 7%. The *CFA U-Net* further improves this metric to 97.61%, an increase of almost 11 percentage points over the Attention U-Net, effectively matching the high-recall performance of the best-in-class standard U-Net (97.40%).

However, compared to the best-in-class MAE of the standard Attention U-Net of 4.14 (with 10-line spacing), the *CFA^S U-Net* registered a significantly higher MAE of 4.77 (an increase of more than 15%). The addition of the Sobel head to *CFA U-Net* mitigated this effect, reducing the MAE to 4.44. Therefore, while the spatial head provides a strong inductive bias for surface completeness, the Sobel head constrains these interpolations geometrically. The *CFA U-Net* represents the most effective

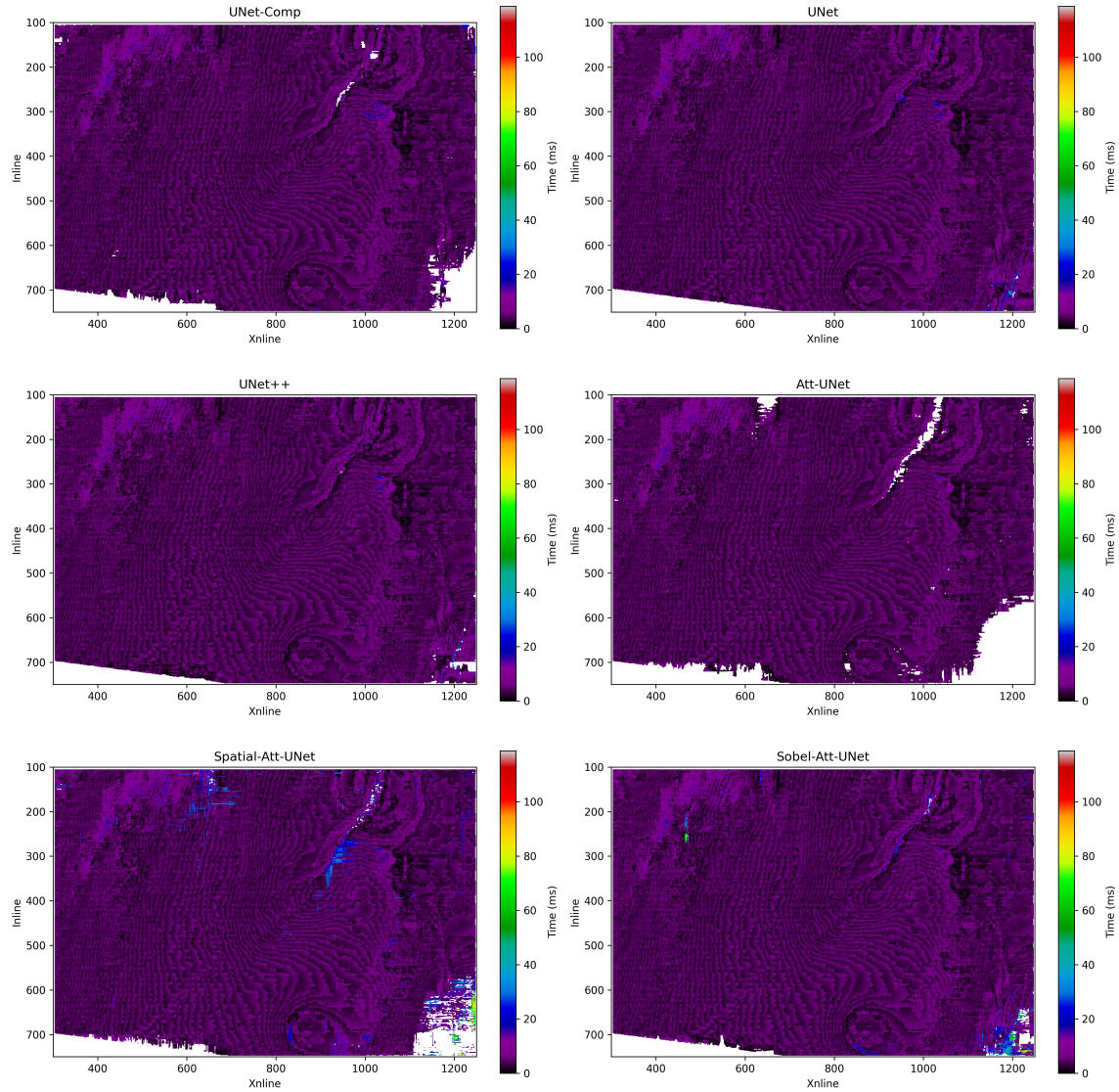


Figure 9: **Spatial distribution of prediction error (TWT difference) for the 10-line spacing scenario.** Each map shows the pixel-wise difference between the predicted and true two-way travel time. Cooler colors indicate minimal error, while warmer colors highlight regions with larger prediction discrepancies. Such visualizations help diagnose model performance, revealing that errors for most architectures are often concentrated along complex fault zones where horizon continuity is disrupted.

trade-off, achieving near-optimal surface coverage while better preserving geometric precision than the spatial-only approach. This approach highlights the crucial need to evaluate both pixel-level and geometric metrics to understand model performance fully.

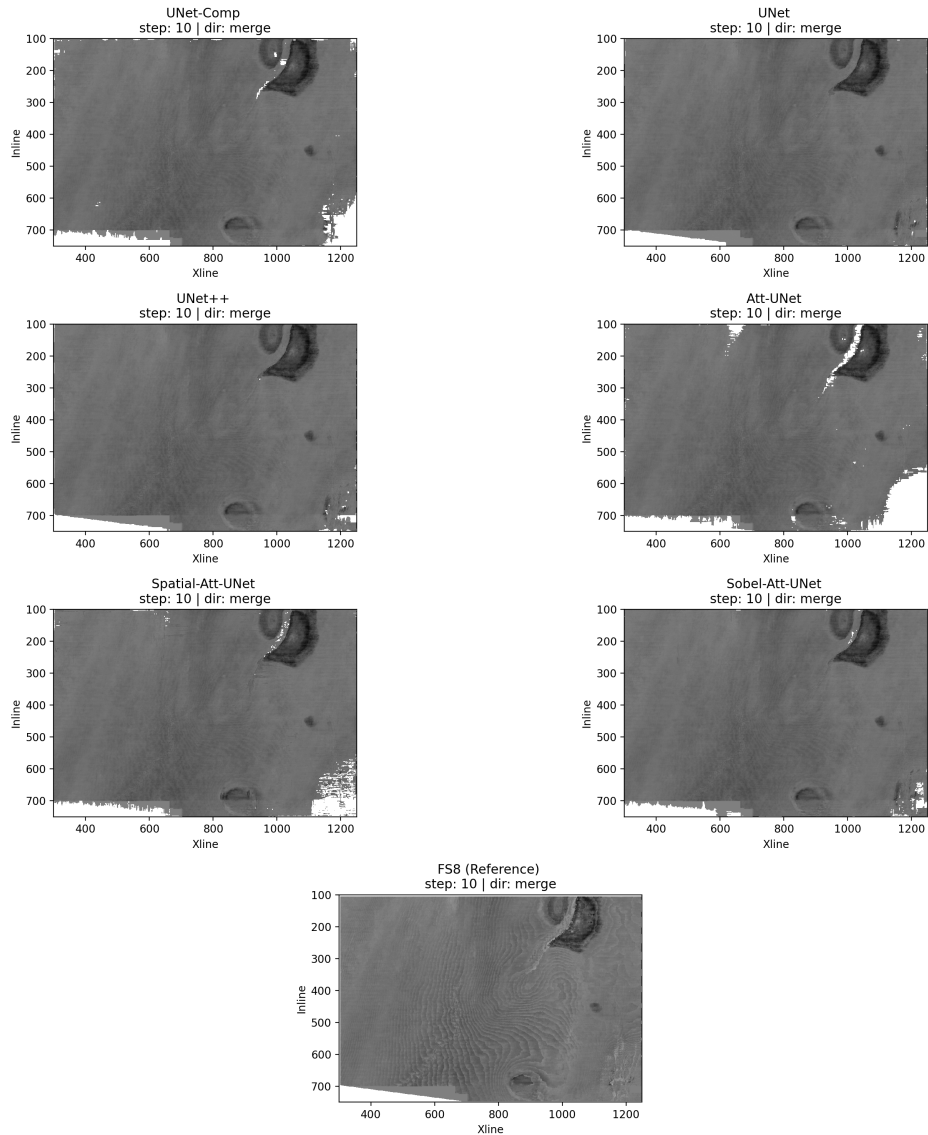


Figure 10: **Map view of predicted horizons on the F3 data set** with merged results for different architectures trained with a 10-line spacing for both inline and crossline directions. The bottom map shows the ground truth (FS8) horizon for reference.

Figure 9 provides a qualitative complement to the metrics in Table 3 by visualizing the spatial distribution of prediction errors (TWT difference) for the 10-line spacing scenario. The error maps reveal the different strategies that the models employ when encountering geological uncertainty. The standard U-Net, with its broad convolutional feature extractors, interpolates across ambiguous regions such as fault zones where reflector continuity is low. This strategy prioritizes surface completeness (high recall), resulting in moderate, localized prediction errors, visualized as the cooler, blue regions representing differences of **10-20 ms**. Figure 10 presents a map

view of the merged horizons for each architecture trained with 10-line spacing in both directions. The observed gaps align with the model coverage reported in Table 3.

In contrast, the Attention U-Net focuses only on high-confidence features. When its attention gate encounters the low-confidence signature of a fault zone, it correctly identifies this region as highly uncertain and actively suppresses the prediction entirely, creating a gap in the output surface. In the error map, these gaps manifest themselves as significant apparent differences (~ 100 ms), not because the model made an inaccurate prediction, but because it correctly identified its own uncertainty and abstained altogether. Therefore, in this case, the Attention U-Net achieves its superior MAE by precisely predicting only what it is sure about, forgoing risky interpolations in complex zones.

The introduction of the spatial-only head in the CFA^S **U-Net** successfully improves the baseline Attention U-Net’s coverage from 93.26% to 96.83% (at 10-line spacing), but this comes with a sharp penalty to precision, as its MAE increases from 4.14 to 4.77. The full CFA **U-Net**, which adds a Sobel filter head, demonstrates a superior solution. It enhances the surface coverage even further to 98.19%, nearly matching the high-recall standard U-Net (98.27%), while the Sobel head constrains the geometry, mitigating the error and yielding a much lower MAE of 4.44. Therefore, the CFA **U-Net** resolves the fundamental trade-off. It achieves the high surface coverage of a standard U-Net while maintaining geometric precision far superior to that of the spatial-only attention model, establishing it as the most balanced and effective architecture for horizon reconstruction in this data set.

Mexilhão Field

To assess architectural robustness, experiments were performed on the Mexilhão Field seismic volume, a data set characterized by significant faulting and lower reflector continuity. The evaluation again relies on comparing both pixel-level accuracy (IoU), which measures classification correctness, and post-processed geometric utility, which measures the real-world precision (MAE) and completeness (surface coverage) of the final horizon.

The increased geological complexity had a clear and quantifiable impact on performance compared to the F3 block. The standard U-Net’s validation IoU at 10-line spacing, for example, dropped by **7.5%** from 0.946 on F3 to 0.875 on Mexilhão (Table 4). This data set also inverted the relationship between geometric errors, as detailed in Table 5. While the Mean Absolute Error (MAE) for the U-Net was **40.2%** lower (decreasing from 4.23 to 2.53), the Mean Squared Error (MSE) was **85.2%** higher (increasing from 23.26 to 43.08). This divergence indicates that while predictions were often closer to the true horizon on average (lower MAE), they were also prone to larger, more significant outlier errors (higher MSE), which can be a direct consequence of the complex geology.

The challenging geology of the Mexilhão Field reveals a sharp trade-off between

Architecture	Spacing	Mean Acc.		Mean IoU		Mean Dice Loss	
		Train	Valid	Train	Valid	Train	Valid
U-Net (Comp.)	10	0.9957	0.9942	0.8543	0.8063	0.4457	0.2848
	20	0.9969	0.9946	0.8809	0.8070	0.3829	0.2954
	40	0.9966	0.9929	0.8729	0.7643	0.4874	0.3675
U-Net	10	0.9993	0.9968	0.9685	0.8747	0.0332	0.1439
	20	0.9995	0.9968	0.9769	0.8724	0.0240	0.1468
	40	0.9987	0.9954	0.9432	0.8240	0.0627	0.2165
U-Net++	10	0.9997	0.9963	0.9877	0.8573	0.0158	0.1690
	20	0.9997	0.9962	0.9866	0.8515	0.0225	0.1791
	40	0.9996	0.9948	0.9825	0.8086	0.1070	0.2813
Attn. U-Net	10	0.9998	0.9963	0.9899	0.8500	0.0122	0.1833
	20	0.9997	0.9959	0.9890	0.8322	0.0385	0.2365
	40	0.9998	0.9944	0.9893	0.7770	0.0642	0.3500
CFA^S U-Net	10	0.9998	0.9963	0.9911	0.8478	0.0154	0.1938
	20	0.9997	0.9955	0.9875	0.8162	0.0459	0.2731
	40	0.9997	0.9944	0.9882	0.7755	0.0780	0.3592
CFA U-Net	10	0.9996	0.9971	0.9827	0.8812	0.0213	0.1404
	20	0.9996	0.9968	0.9841	0.8695	0.0195	0.1554
	40	0.9995	0.9959	0.9789	0.8386	0.0449	0.2201

Table 4: **Segmentation metrics on the Mexilhão data set.** "Valid" refers to the validation set. Values represent the mean of models trained on inline and crossline data. U-Net (Comp.) and Attn. U-Net represents Compressed U-Net and Attention U-Net, respectively.

geometric precision and surface completeness. The results, summarized in Table 5, show that no single architecture dominates on both metrics. Instead, a clear hierarchy emerges where the proposed CFA U-Net excels at precision, while the U-Net++ architecture is the unambiguous leader in recall.

On this challenging data set, the proposed CFA U-Net emerges as the **top-performing architecture in terms of overall precision** when the data is reasonably dense. At the 10-line spacing, it achieved both the highest validation IoU (0.881) and the lowest MAE (2.49), demonstrating superior pixel-level and geometric accuracy (Tables 4 and 5). This MAE is **1.6%** lower than the standard U-Net’s (2.53) and **6.7%** lower than the U-Net++’s (2.67). This result is significant, as it shows that the multi-head attention mechanism, which provides spatial and edge-aware inductive biases, is highly effective at navigating the complex faulting that challenges other models.

While the CFA U-Net was the most precise, the U-Net++ **architecture proved to be the best approach for recall**, consistently delivering the most complete surfaces, as quantified in Table 5 and visualized in Figure 11. At 10-line spacing, it

Architecture	Spacing	MAE	MSE	Area (%)
U-Net (Comp.)	10	2.89	60.79	94.43
	20	2.85	61.32	91.03
	40	3.93	126.51	91.17
U-Net	10	2.53	43.08	95.86
	20	2.64	48.52	95.41
	40	3.11	72.14	94.27
U-Net++	10	2.67	52.87	96.25
	20	2.93	66.53	95.82
	40	4.17	146.23	94.79
Attn. U-Net	10	2.73	59.10	90.35
	20	3.21	93.55	87.30
	40	3.56	119.03	78.70
<i>CFA</i> ^S U-Net	10	2.67	55.37	89.96
	20	3.48	117.55	85.15
	40	3.32	98.92	79.45
<i>CFA</i> U-Net	10	2.49	45.30	92.97
	20	2.94	72.18	92.12
	40	3.42	106.97	89.34

Table 5: Geometric errors (MAE, MSE) and surface coverage area of the merged and filtered results for each model on the Mexilhão data set, including the proposed *CFA*^S U-Net with spatial head and *CFA* U-Net including semantic, spatial, and sobel heads.

achieved a surface coverage of **96.25%**, the highest of any model, which is a notable improvement over the standard U-Net (95.86%) and significantly higher than the high-precision *CFA* U-Net (92.97%). However, this superior recall came at the cost of precision, with an MAE of 2.67 that was higher than both the standard U-Net and the *CFA* U-Net. In contrast, the standard **Attention U-Net produced** one of the least complete surfaces (coverage of 90.35%) without achieving a competitive level of geometric precision (MAE of 2.73), highlighting its limitations in regions of widespread geological ambiguity.

A more detailed, cross-sectional view of these predictions is provided in Figure 12. This 2D profile view allows for a direct comparison of the predicted horizon (red line) from each architecture against the ground truth (dotted yellow line). These profiles visually corroborate the quantitative metrics, illustrating how the predictions from models like U-Net++, the standard U-Net, Attention U-Net, and both Context-Fusion Attention U-Net closely follow the true horizon across large continuous segments, while also revealing the small-scale geometric deviations that contribute to their respective MAE scores.

The performance hierarchy changes significantly under the most challenging condi-

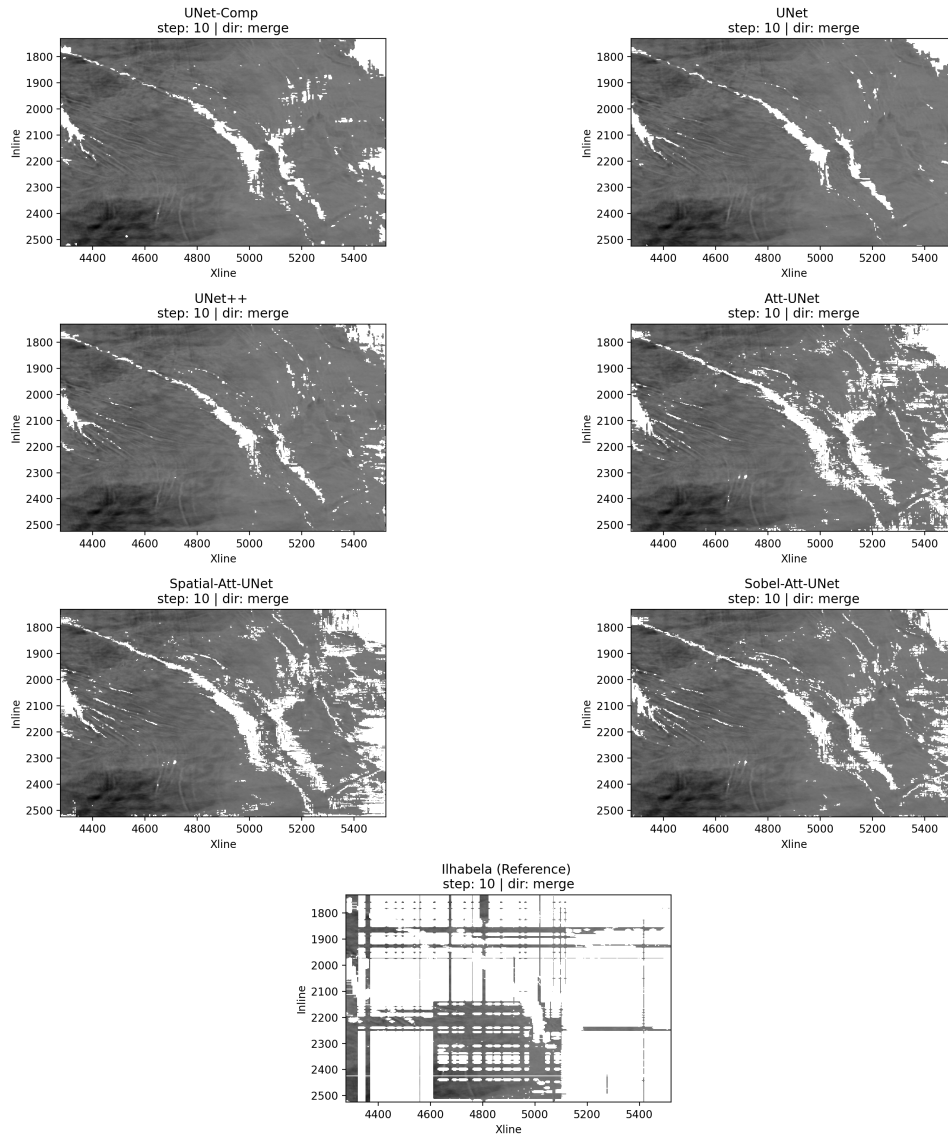


Figure 11: **Map view of predicted horizons on the Mexilhão data set.** The panel shows the merged results for each architecture trained with a 10-line spacing for both directions (the inline and cross-line), as indicated by its title. Models (a) Compressed U-Net, (b) U-Net, (c) U-Net++, (d) Attention U-Net, (e) Spatial CFA U-Net, (f) Context-Fusion Attention (Sobel, Spatial and Context) U-Net, and (g) the ground truth interpretation, highlighting the irregular grid used for training samples.

tion of extreme data sparsity (40-line spacing), as detailed in Table 5. As the training set is reduced to every 40th line, the architectural complexity that benefited *CFA* U-Net becomes a liability. Its geometric precision declines, with its MAE increasing to 3.42. In this low-data regime, the simpler **standard U-Net demonstrates superior robustness, reclaiming the top spot for geometric precision with the lowest MAE of 3.11.** This result is nearly **10%** better than the *CFA* U-Net’s MAE,

suggesting that the standard U-Net’s less complex convolutional structure provides a more substantial, more generalizable inductive bias when training data is severely limited. Even in this sparse scenario, U-Net++ maintains its status as the recall leader with a coverage of 94.79%.

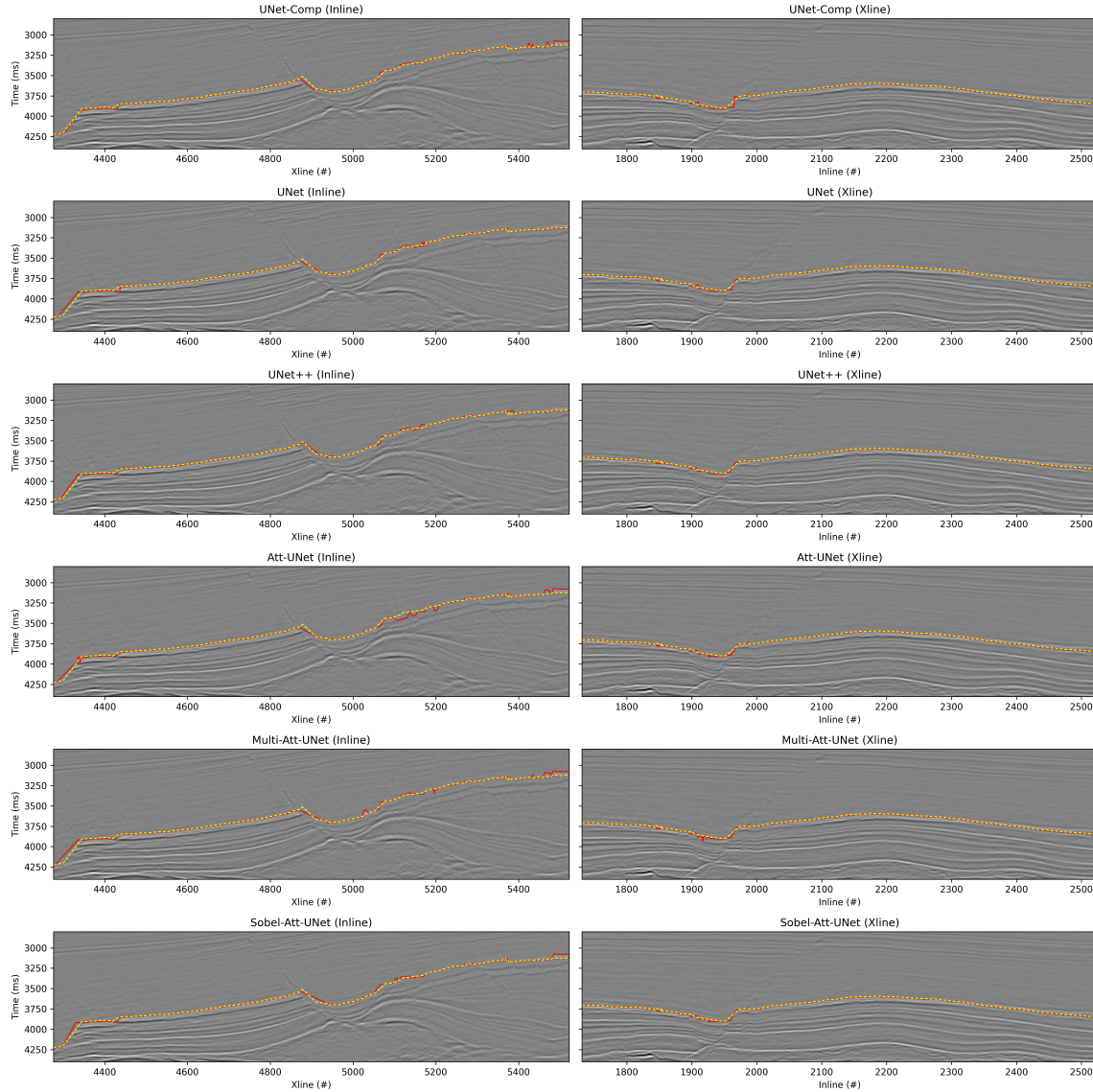


Figure 12: **2D profile comparison of predicted horizons.** This figure shows a cross-section of the results on a specific inline (left) and crossline (right). The predicted horizon (red line) from each architecture is compared against the true interpreted horizon (dotted yellow line).

In conclusion, while the sophisticated attention mechanism of the *CFA U-Net* makes it the best model for navigating complex geology with adequate data, the architectural simplicity of the **standard U-Net** makes it more robust and reliable under conditions of extreme data scarcity. The **U-Net++** remains the most practical choice in any scenario where maximizing surface continuity was the primary objective.

This nuanced, task-dependent hierarchy highlights the critical importance of selecting an architecture based not only on the geological setting but also on the density of the available training data.

Discussion

The U-Net-based architectures have been systematically evaluated, confirming that the optimal choice for seismic interpretation is highly dependent on both geological complexity and data sparsity. The evidence is found in the performance drop between the simple F3 block and the complex Mexilhão field, where the peak validation IoU for the best-performing model fell from 0.946 (Table 2) to 0.881 (Table 4). The data scarcity further compounded this challenge, as seen on the complex Mexilhão data where the standard U-Net’s geometric error (MAE) increased by nearly 23% when the training data was reduced (Table 5).

The analysis of the established architectures revealed distinct, specialized roles. The standard U-Net proved to be a robust all-rounder, while U-Net++ consistently delivered the highest surface coverage (e.g., 96.3% on Mexilhão, Table 5), making it the best choice for large-scale structural mapping. The Attention U-Net confirmed its status as a high-precision specialist on simple data, achieving the lowest geometric error (MAE of 4.14 ms) on the F3 data set (Table 3). However, its cautious prediction strategy, which relies on high-confidence features, proved to be a critical weakness in complex geology, a trade-off visually apparent in the surface gaps shown in Figure 8. On the Mexilhão data set, its surface coverage fell to a low of 78.7% (Table 5) as the attention gates learned to suppress predictions in faulted zones.

The design of the proposed Context-Fusion Attention (CFA) U-Net was validated through an ablation study. An intermediate model with only a spatial head (CFA^S U-Net) successfully addressed the low-recall problem of the standard Attention U-Net, increasing coverage from 93.3% to 96.8% at 10-line spacing. However, this came at the cost of a sharp increase 15% in geometric error and the highest MAE scores, which reached a peak at 5.18 ms (Table 3). This finding proved that adding a spatial prior alone is insufficient, thus validating the need for the Sobel head in the full CFA U-Net to create a more balanced architecture.

The success of the full CFA U-Net was first demonstrated by its class-leading robustness in the F3 data set. It completely resolved the recall issue of the baseline attention model, boosting its surface coverage at 40-line spacing by a massive 10.8 percentage points. This improvement was so significant that it achieved the highest surface coverage (97.6%) of all architectures in this data-sparse scenario, outperforming even the high-recall of U-Net++ (97.3%), as shown in Table 3 and visualized in Figure 8. Therefore, the context-aware design provides a superior inductive bias, allowing the model to generalize exceptionally well from limited data.

The performance of the CFA U-Net was fully assessed on the challenging Mexilhão data set. The CFA U-Net enhanced precision, providing both the highest IoU valida-

tion (0.881, Table 4) and the lowest MAE (2.49 ms, Table 5) of any model. This dual achievement, visually corroborated by the map views and 2D profiles in Figures 11 and 12, confirms that its hybrid attention mechanism is uniquely effective at navigating geologically complex regions. Therefore, the validation of the CFA U-Net and the hybrid workflow shows that context-aware attention mechanisms are critical components to creating more robust, accurate, and practical tools for automated seismic interpretation.

Conclusion

The results established that while standard models like U-Net and U-Net++ are robust baselines that excel at all-around performance and high-recall mapping, respectively, the standard Attention U-Net is a high-precision specialist whose practical utility is limited by poor surface completeness on complex geology. This precision-recall trade-off motivated the development of our proposed context-aware architecture.

The CFA U-Net successfully addressed these limitations, demonstrating class-leading robustness on the F3 data set. The model improved the critical recall issue of the baseline attention model, boosting surface coverage at 40-line spacing by a massive 10.8 percentage points. Furthermore, the improvement was so significant that it achieved the highest surface coverage (97.6%) of all architectures in this data-sparse scenario, outperforming even the high-recall champion, U-Net++ (97.3%). Therefore, the context-aware design provides a superior inductive bias for generalizing from limited data.

Furthermore, the CFA U-Net’s capabilities were also assessed on the challenging Mexilhão data set. The results demonstrate high precision, yielding both the highest validation IoU (0.881) and the lowest MAE (2.49 ms) among all models. Therefore, the hybrid fusion attention mechanism is uniquely effective at navigating geologically complex regions to produce the most accurate and reliable interpretations as compared to the Attention U-Net.

Ultimately, this robust hybrid workflow successfully combines neural networks with DBSCAN and merged orthogonal (inline and crossline) predictions for comprehensive 3D horizon reconstruction. The development of context-aware attention mechanisms was validated through the proposed CFA U-Net, which is considered a promising path toward creating practical interpretation tools.

DATA AND MATERIALS AVAILABILITY

Data associated with this research are available and can be obtained by contacting the corresponding author.

CORRESPONDING AUTHOR

Correspondence and requests for materials should be addressed to *Dr. Jose Luis Silva* at `jse Luis.silva@gmail.com`.

ACKNOWLEDGEMENTS

This work was supported by the National Council for Scientific and Technological Development (CNPq), Brazil, under grants No. 409718/2022-0 and No. 445344/2024-5.

REFERENCES

- AlSalmi, H., and A. H. Elsheikh, 2024, Automated seismic semantic segmentation using attention u-net: *Geophysics*, **89**, WA247–WA263.
- Barbosa, M. R. S., V. Carneiro, and A. G. Cerqueira, 2024, Direct hydrocarbon indicators mapping via joint cluster analysis: A two-step approach over 3d seismic data: *Revista de Geociências do Nordeste (Northeast Geosciences Journal)*, **10**.
- Bergmeir, C., and J. M. Benítez, 2012, On the use of cross-validation for time series predictor evaluation: *Information Sciences*, **191**, 192–213.
- Calhes, D., F. K. Kobayashi, A. B. Mattos, M. M. Macedo, and D. A. Oliveira, 2021, Simplifying horizon picking using single-class semantic segmentation networks: 2021 34th SIBGRAPI Conference on Graphics, Patterns and Images (SIBGRAPI), IEEE, 286–292.
- Davis, J., and M. Goadrich, 2006, The relationship between precision-recall and roc curves: *Proceedings of the 23rd international conference on Machine learning*, 233–240.
- Deng, D., 2020, Dbscan clustering algorithm based on density: 2020 7th international forum on electrical engineering and automation (IFEEA), IEEE, 949–953.
- Dhua, A., D. N. Sarma, S. Singh, and B. Roy, 2015, Segmentation of images using density-based algorithms: *International Journal of Advanced Research in Computer and Communication Engineering*, **4**, 273–277.
- Dorn, G. A., 1998, Modern 3-d seismic interpretation: *The Leading Edge*, **17**, 1262–1262.
- Ester, M., H.-P. Kriegel, J. Sander, X. Xu, et al., 1996, A density-based algorithm for discovering clusters in large spatial databases with noise: *kdd*, 226–231.
- Hale, D., 2013, Dynamic warping of seismic images: *GEOPHYSICS*, **78**, S105–S115.
- Harrigan, E., J. Kroh, W. Sandham, and T. Durrani, 1992, Seismic horizon picking using an artificial neural network: [Proceedings] ICASSP-92: 1992 IEEE International Conference on Acoustics, Speech, and Signal Processing, IEEE, 105–108 vol.3.
- Ioffe, S., and C. Szegedy, 2015, Batch normalization: Accelerating deep network training by reducing internal covariate shift: *International conference on machine learning*, pmlr, 448–456.
- Jadon, S., 2020, A survey of loss functions for semantic segmentation: *arXiv preprint arXiv:2006.14822*.
- Kaufman, S., S. Rosset, C. Perlich, and O. Stitelman, 2012, Leakage in data mining: Formulation, detection, and avoidance: *ACM Transactions on Knowledge Discovery from Data (TKDD)*, **6**, 1–21.
- Kingma, D. P., and J. Ba, 2014, Adam: A method for stochastic optimization: *arXiv preprint arXiv:1412.6980*.
- Krizhevsky, A., I. Sutskever, and G. E. Hinton, 2012, Imagenet classification with deep convolutional neural networks: *Advances in neural information processing systems*, **25**.
- Kusuma, T., and B. C. Fish, 1993, Toward more robust neural-network first break and horizon pickers: *SEG Technical Program Expanded Abstracts 1993*, Society of

- Exploration Geophysicists, 238–241.
- LeCun, Y., Y. Bengio, and G. Hinton, 2015, Deep learning: *Nature*, **521**, 436–444.
- Luo, Y., G. Zhang, J. Zhang, Y. Li, Y. Lin, B. Li, C. Liang, and L. Li, 2023, Sequence-constrained multitask horizon tracking: *GEOPHYSICS*, **88**, IM15–IM27.
- Marfurt, K. J., V. Sudhaker, A. Gersztenkorn, K. D. Crawford, and S. E. Nissen, 1999, Coherency calculations in the presence of structural dip: *GEOPHYSICS*, **64**, 104–111.
- Mattos, A. B., D. Civitarese, D. Szwarcman, M. Oliveira, S. Zaytsev, D. G. Semin, and D. A. B. Oliveira, 2021, Enabling robust horizon picking from small training sets: *IEEE Transactions on Geoscience and Remote Sensing*, **59**, 5317–5324.
- Nair, V., and G. E. Hinton, 2010, Rectified linear units improve restricted boltzmann machines: *Proceedings of the 27th international conference on machine learning (ICML-10)*, 807–814.
- Oktay, O., J. Schlemper, L. L. Folgoc, M. Lee, M. Heinrich, K. Misawa, K. Mori, S. McDonagh, N. Y. Hammerla, B. Kainz, et al., 2018, Attention u-net: Learning where to look for the pancreas: *arXiv preprint arXiv:1804.03999*.
- Peters, B., J. Granek, and E. Haber, 2019, Multiresolution neural networks for tracking seismic horizons from few training images: *Interpretation*, **7**, SE201–SE213.
- Poulinakis, K., D. Drikakis, I. W. Kokkinakis, and S. M. Spottswood, 2023, Machine-learning methods on noisy and sparse data: *Mathematics*, **11**, 236.
- Ravasi, M., and C. Birnie, 2022, A joint inversion-segmentation approach to assisted seismic interpretation: *Geophysical Journal International*, **228**, 893–912.
- Roberts, D. R., V. Bahn, S. Ciuti, M. S. Boyce, J. Elith, G. Guillera-Arroita, S. Hauenstein, J. J. Lahoz-Monfort, B. Schröder, W. Thuiller, et al., 2017, Cross-validation strategies for data with temporal, spatial, hierarchical, or phylogenetic structure: *Ecography*, **40**, 913–929.
- Ronneberger, O., P. Fischer, and T. Brox, 2015, U-net: Convolutional networks for biomedical image segmentation: *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5–9, 2015, proceedings, part III 18*, Springer, 234–241.
- Schroot, B. M., G. T. Klaver, and R. T. Schüttenhelm, 2005, Surface and subsurface expressions of gas seepage to the seabed—examples from the southern north sea: *Marine and Petroleum Geology*, **22**, 499–515.
- Schubert, E., J. Sander, M. Ester, H. P. Kriegel, and X. Xu, 2017, Dbscan revisited, revisited: why and how you should (still) use dbscan: *ACM Transactions on Database Systems (TODS)*, **42**, 1–21.
- Sheriff, R. E., and L. P. Geldart, 1995, *Exploration seismology*: Cambridge university press.
- Stark, T. J., 2003, Unwrapping instantaneous phase to generate a relative geologic time volume: *SEG Technical Program Expanded Abstracts 2003*, Society of Exploration Geophysicists, 1707–1710.
- Stollenga, M. F., W. Byeon, M. Liwicki, and J. Schmidhuber, 2015, Parallel multi-dimensional lstm, with application to fast biomedical volumetric image segmentation: *Advances in neural information processing systems*, **28**.
- Troccoli, E. B., A. G. Cerqueira, J. B. Lemos, and M. Holz, 2022, K-means cluster-

- ing using principal component analysis to automate label organization in multi-attribute seismic facies analysis: *Journal of Applied Geophysics*, **198**, 104555.
- Tschannen, V., M. Delescluse, N. Ettrich, and J. Keuper, 2020, Extracting horizon surfaces from 3d seismic data using deep learning: *GEOPHYSICS*, **85**, N17–N26.
- Wu, H., and B. Zhang, 2018, A deep convolutional encoder-decoder neural network in assisting seismic horizon tracking: arXiv preprint arXiv:1804.06814.
- Wu, H., B. Zhang, T. Lin, D. Cao, and Y. Lou, 2019, Semiautomated seismic horizon interpretation using the encoder-decoder convolutional neural network: *GEOPHYSICS*, **84**, B403–B417.
- Wu, X., and S. Fomel, 2018, Least-squares horizons with local slopes and multigrid correlations: *GEOPHYSICS*, **83**, IM29–IM40.
- Wu, X., and G. Zhong, 2012, Generating a relative geologic time volume by 3d graph-cut phase unwrapping method with horizon and unconformity constraints: *GEOPHYSICS*, **77**, O21–O34.
- Yang, L., and S. Z. Sun, 2020, Seismic horizon tracking using a deep convolutional neural network: *Journal of Petroleum Science and Engineering*, **187**, 106709.
- Yu, S., and J. Ma, 2021, Deep learning for geophysics: Current and future trends: *Reviews of Geophysics*, **59**, e2021RG000742.
- Zeiler, M. D., and R. Fergus, 2014, Visualizing and understanding convolutional networks: *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part I 13*, Springer, 818–833.
- Zhou, Z., M. Mahfuzur, R. Siddiquee, N. Tajbakhsh, , and J. Liang, 2018, Unet++: A nested u-net architecture for medical image segmentation: arXiv:1807.10165.