

Cloud Base Height Retrieval: Sprints 3 to 5 Status Report

Rylan Malarchick

November 2025

Abstract

This report documents the completion of Feature Engineering & Integration, along with Hybrid Model Development, as well as Part 5 (Advanced Deep Learning) for the Cloud Base Height (CBH) retrieval project. It presents a comprehensive evaluation spanning classical machine learning baselines, hybrid CNN architectures, and state-of-the-art deep learning models with temporal modeling and advanced fusion techniques. The key finding is that **temporal Vision Transformer (ViT) models with physics-informed consistency losses are the first deep learning architectures to outperform the physical baseline**, achieving $R^2 = 0.728$ with 126-meter mean absolute error. All results reported use real operational data from NASA ER-2 flights with validated ERA5 atmospheric reanalysis.

Contents

1 Summary	3
1.1 Overview	3
1.2 Key Results Summary	3
1.3 Breakthrough Achievement	3
1.4 Critical Findings	3
2 Section 3: Feature Engineering & Integration	5
2.1 Part 1: Geometric Features	5
2.1.1 Shadow-Based CBH Derivation	5
2.2 Part 2: Atmospheric Features	6
2.2.1 ERA5 Reanalysis Integration	6
2.3 Part 3: Physical Baseline Model	7
2.3.1 Model Architecture	7
2.3.2 Validation Protocol	7
2.3.3 Results	8
2.4 Integrated Feature Store	8
2.5 Feature Importance Analysis	9
3 Section 4: Hybrid Model Development	10
3.1 Motivation	10
3.2 Part 1: CNN Architecture Development	10
3.2.1 Model 1: Image-Only Baseline CNN	10
3.2.2 Model 2: Concatenation Fusion	11
3.2.3 Model 3: Attention Fusion	11

3.3 Ablation Study	12
4 Section 5: Advanced Deep Learning	13
4.1 Overview	13
4.2 Part 1: Pre-Trained Backbones	13
4.2.1 Model: ResNet-50 Baseline	13
4.2.2 Model: ViT-Tiny Baseline	14
4.3 Part 2: Temporal Modeling	15
4.3.1 Multi-Frame Temporal ViT	15
4.3.2 Physics-Informed Temporal Consistency Loss	16
4.4 Part3: Advanced Fusion (FiLM)	17
4.4.1 Feature-wise Linear Modulation	17
4.5 Part 5 Summary	19
5 Cross-Part/Task Analysis	20
5.1 Performance Evolution	20
5.2 Feature Importance: Physical vs. Learned	20
5.3 Computational Requirements	21
5.4 Validation Protocol Robustness	21
6 Data	22
6.1 ERA5 Processing Details	22
7 Deployment Recommendations	23
7.1 Recommended Model Architecture	23
7.2 Alternative: Dual-Model Deployment	23
7.3 Known Limitations and Mitigations	24
8 Future Work	25
8.1 High Priority	25
8.2 Medium Priority (Research Extensions)	25
8.3 Low Priority (Exploratory Research)	26
9 Conclusions	27
9.1 Summary of Achievements	27
9.2 Critical Findings	27
9.3 Final Remarks	27

1 Summary

1.1 Overview

- **Section 3:** Integration of geometric features (WP1) with atmospheric features (WP2), establishing physical baseline and validation framework
- **Section 4:** Development and evaluation of hybrid deep learning models combining image features with physical constraints using attention-based fusion
- **Section 5:** Advanced deep learning with pre-trained backbones, temporal modeling, physics-informed losses, and multimodal fusion (FiLM)

1.2 Key Results Summary

Performance Evolution Across Sections:

Model	Section	R^2	MAE (m)
Physical GBDT (baseline)	3	0.668	137
Image-only CNN	4	0.279	233
Attention Fusion CNN	4	0.326	221
ResNet-50 Baseline	5	0.524	171
ViT-Tiny Baseline	5	0.577	166
Temporal ViT	5	0.727	126
Temporal ViT + Consistency ($\lambda = 0.1$)	5	0.728	126
FiLM Fusion	5	0.542	166

1.3 Breakthrough Achievement

First Deep Learning Model to Beat Physical Baseline:

The Temporal ViT model with physics-informed temporal consistency loss represents a breakthrough:

- **Performance:** $R^2 = 0.728$ (vs. 0.668 baseline) = +9% improvement
- **Accuracy:** MAE = 126 m (vs. 137 m baseline) = 11 m improvement (8%)
- **Architecture:** ViT-Tiny encoder with 5-frame temporal attention
- **Innovation:** Physics-informed temporal consistency regularization

1.4 Critical Findings

1. **Temporal information is essential:** Single-frame models ($R^2 \approx 0.28\text{--}0.58$) cannot compete with temporal models ($R^2 = 0.73$)
2. **Pre-trained backbones improve over CNNs from scratch:** ViT-Tiny (ImageNet-21k pre-trained) achieves $R^2 = 0.577$ vs. custom CNN $R^2 = 0.279$
3. **Attention mechanisms outperform CNNs:** Vision Transformers systematically outperform ResNet architectures

4. **Physics-informed losses provide marginal gains:** Temporal consistency loss ($\lambda = 0.1$) improves R^2 from 0.727 to 0.728
5. **ERA5 fusion remains challenging:** FiLM fusion ($R^2 = 0.542$) underperforms image-only ViT, suggesting atmospheric features require better integration strategies

2 Section 3: Feature Engineering & Integration

2.1 Part 1: Geometric Features

2.1.1 Shadow-Based CBH Derivation

Building on solar geometry analysis, I implemented shadow-length-based cloud base height estimation following the physical principle:

$$H_{\text{cloud}} = L_{\text{shadow}} \cdot \tan(\theta_{\text{SZA}}) \quad (1)$$

where L_{shadow} is the detected shadow length (meters) and θ_{SZA} is the solar zenith angle.

Implementation:

- Edge detection on 440×640 grayscale images using Canny algorithm
- Shadow region identification via brightness thresholding (adaptive percentile-based)
- Length measurement in pixel coordinates with camera geometry correction
- Conversion to physical units using aircraft altitude and viewing angle
- Confidence scoring based on edge sharpness and detection consistency

Feature set (10 features):

1. `derived_geometric_H`: Shadow-derived cloud base height estimate
2. `shadow_length_pixels`: Raw detected shadow length
3. `shadow_detection_confidence`: Quality score [0, 1]
4. `sza_rad`: Solar zenith angle (radians)
5. `saa_rad`: Solar azimuth angle (radians)
6. `cloud_top_edge_y`: Vertical position of cloud top edge
7. `cloud_bottom_edge_y`: Vertical position of cloud bottom edge
8. `edge_sharpness`: Mean gradient magnitude along edges
9. `altitude_m`: Aircraft altitude (from GPS)
10. `geometric_consistency`: Multi-frame consistency metric

Data quality:

- 87.1% of samples have valid shadow detections ($\text{confidence} > 0.5$)
- 12.9% missing values handled via median imputation
- Failures occur in: (1) optically thin cirrus, (2) broken cloud fields, (3) $\text{SZA} > 70$ (this is a marginal part of our dataset so not a huge concern - and also does not really correspond to any of our CPL picks)

2.2 Part 2: Atmospheric Features

2.2.1 ERA5 Reanalysis Integration

Data Source:

ERA5 atmospheric reanalysis data was processed from NASA archives:

- **Surface-level data:** 119 daily files (Oct 23, 2024 – Feb 19, 2025)
- **Pressure-level data:** 37 vertical levels, hourly temporal resolution
- **Spatial resolution:** 0.25×0.25 (approximately 25 km)
- **Coverage:** All 5 flight dates fully covered
- **Processing success:** 933/933 samples (100%)

Spatiotemporal Matching:

- **Spatial:** Nearest neighbor interpolation to flight track GPS coordinates
- **Temporal:** Nearest hourly ERA5 timestamp to image acquisition time
- **Vertical:** Boundary layer features extracted from surface and pressure levels

Derived atmospheric features (9 features):

1. `blh_m`: Boundary layer height (from ERA5 BLH field)
2. `lcl_m`: Lifting condensation level (computed from T, Td)
3. `inversion_height_m`: Temperature inversion base height
4. `moisture_gradient`: Vertical moisture gradient (kg/kg/m)
5. `stability_index`: Atmospheric stability (lapse rate, K/km)
6. `surface_temp_k`: 2-meter temperature
7. `surface_dewpoint_k`: 2-meter dewpoint temperature
8. `surface_pressure_pa`: Surface pressure
9. `profile_confidence`: ERA5 data quality indicator

Feature Statistics (Real ERA5):

Physical validation:

- BLH values (658 ± 485 m) consistent with marine boundary layer
- Stability index (3.81 K/km) indicates stable atmosphere (less than standard 6.5 K/km)
- LCL correlates with observed low cloud base heights (mean CBH = 830 m)
- All values physically realistic for coastal/oceanic flight conditions

Variable	Mean	Std Dev
Boundary Layer Height	658 m	485 m
Lifting Condensation Level	839 m	589 m
Inversion Height	875 m	688 m
Moisture Gradient	-1.07×10^{-6} kg/kg/m	—
Stability Index	3.81 K/km	0.93 K/km
Surface Temperature	284.4 K	9.1 K
Surface Dewpoint	277.7 K	9.5 K
Surface Pressure	96,928 Pa	7,535 Pa

2.3 Part 3: Physical Baseline Model

2.3.1 Model Architecture

Algorithm: XGBoost Gradient Boosted Decision Trees (GBDT)

Input features: 19 total

- 10 geometric features (WP1)
- 9 atmospheric features (WP2)

Hyperparameters:

- `max_depth`: 6 (prevents overfitting on small dataset)
- `learning_rate`: 0.1 (conservative step size)
- `n_estimators`: 100 (ensemble size)
- `subsample`: 0.8 (row sampling for regularization)
- `colsample_bytree`: 0.8 (column sampling)
- `min_child_weight`: 3 (minimum samples per leaf)
- `gamma`: 0.1 (minimum loss reduction for split)
- `reg_alpha`: 0.01 (L1 regularization)
- `reg_lambda`: 1.0 (L2 regularization)

2.3.2 Validation Protocol

Method: Stratified 5-Fold Cross-Validation

Rationale:

- Ensures balanced CBH distribution across folds (10 quantile bins)
- Prevents extreme domain shift (Flight F4 has mean CBH = 0.697 km vs. 0.917 km for F0)
- Provides stable performance estimates for hyperparameter tuning
- Leave-One-Flight-Out CV explicitly avoided after catastrophic failure ($R^2 = -3.13$ on F4)

2.3.3 Results

Aggregate Performance:

- Mean R^2 : 0.6759 ± 0.0442
- Mean MAE: 0.1356 ± 0.0068 km (136 meters)
- Mean RMSE: 0.2105 ± 0.0123 km (211 meters)

Per-Fold Results:

Fold	N_{test}	R^2	MAE (km)
0	187	0.629	0.139
1	187	0.753	0.126
2	187	0.663	0.144
3	186	0.641	0.140
4	186	0.693	0.129

Interpretation:

- Low cross-fold variance indicates robust generalization
- Best fold (Fold 1): $R^2 = 0.753$ suggests upper performance bound
- Worst fold (Fold 0): $R^2 = 0.629$ still exceeds earlier hybrid models
- This model establishes the baseline to beat: $R^2 = 0.668$

2.4 Integrated Feature Store

File: sow_outputs/integrated_features/Integrated_Features.hdf5

Structure:

```
Integrated_Features.hdf5
|-- geometric_features/ [933 x 10]
|   |-- derived_geometric_H
|   |-- shadow_length_pixels
|   +-- ... (8 more)
|-- atmospheric_features/ [933 x 9]
|   |-- blh_m
|   |-- lcl_m
|   +-- ... (7 more)
|-- metadata/
|   |-- sample_id [933]
|   |-- flight_id [933]
|   |-- cbh_km (target) [933]
|   |-- latitude [933]
|   |-- longitude [933]
|   +-- timestamp [933]
-- image_features/ [reserved for CNN embeddings]
```

Dataset Statistics:

- **Total samples:** 933 (post-filtering for quality)
- **CBH range:** [0.120, 1.950] km
- **CBH mean:** 0.830 ± 0.371 km
- **Flight distribution:**
 - F0 (30Oct24): 501 samples (mean CBH = 0.917 km)
 - F1 (10Feb25): 191 samples (mean CBH = 0.695 km)
 - F2 (23Oct24): 105 samples (mean CBH = 0.503 km)
 - F3 (12Feb25): 92 samples (mean CBH = 1.081 km)
 - F4 (18Feb25): 44 samples (mean CBH = 0.697 km)

2.5 Feature Importance Analysis

Top 10 Features (GBDT Gain Importance):

Rank	Feature	Type	Importance
1	blh_m	Atmospheric	0.243
2	derived_geometric_H	Geometric	0.187
3	sza_rad	Geometric	0.124
4	lcl_m	Atmospheric	0.098
5	shadow_length_pixels	Geometric	0.076
6	surface_temp_k	Atmospheric	0.061
7	edge_sharpness	Geometric	0.054
8	stability_index	Atmospheric	0.047
9	cloud_bottom_edge_y	Geometric	0.039
10	inversion_height_m	Atmospheric	0.031

Key insights:

- BLH (boundary layer height) is the strongest predictor, validating atmospheric state importance
- Shadow-derived geometric height is second-most important
- Solar zenith angle critical (affects shadow geometry)
- Mix of geometric and atmospheric features in top 10 validates multi-modal approach

3 Section 4: Hybrid Model Development

3.1 Motivation

Section 3 demonstrated that physical features (geometric + atmospheric) achieve strong performance ($R^2 = 0.676$). However, the question remained: *Can deep learning on raw images match or exceed this performance?*

Section 4 addresses this by developing hybrid CNN architectures that combine:

- **Image features:** Learned representations from 440×640 grayscale images
- **Physical features:** Hand-engineered geometric and atmospheric features

3.2 Part 1: CNN Architecture Development

3.2.1 Model 1: Image-Only Baseline CNN

Architecture:

- **Input:** $1 \times 440 \times 640$ grayscale images
- **Encoder:** 4-stage 2D CNN
 - Stage 1: Conv2d($1 \rightarrow 64$, k=3), BatchNorm, ReLU, MaxPool(2)
 - Stage 2: Conv2d($64 \rightarrow 128$, k=3), BatchNorm, ReLU, MaxPool(2)
 - Stage 3: Conv2d($128 \rightarrow 256$, k=3), BatchNorm, ReLU, MaxPool(2)
 - Stage 4: Conv2d($256 \rightarrow 256$, k=3), BatchNorm, ReLU, AdaptiveAvgPool(1, 1)
- **Embedding:** 256-dimensional image representation
- **Regressor:** Linear($256 \rightarrow 128$), ReLU, Dropout(0.3), Linear($128 \rightarrow 1$)
- **Output:** CBH prediction (scalar, km)

Training:

- Optimizer: Adam (lr=1e-4)
- Loss: MSE
- Batch size: 16
- Early stopping: patience=10 epochs
- Data augmentation: Random horizontal flip, brightness jitter ($\pm 10\%$)

Results:

- $R^2: 0.279 \pm 0.031$
- MAE: 0.233 ± 0.019 km (233 meters)
- RMSE: 0.315 ± 0.017 km

Analysis:

- Significantly underperforms physical baseline (R^2 gap = 0.40)
- CNN struggles to learn geometric relationships from scratch
- Limited training data (933 samples) insufficient for CNN generalization
- No pre-training or transfer learning applied

3.2.2 Model 2: Concatenation Fusion

Architecture:

- Same CNN encoder as Model 1 (256-dim image embedding)
- Physical features: 19-dim vector (geometric + atmospheric)
- **Fusion:** Simple concatenation: [img_emb; phys_feat] → 275-dim
- Regressor: Linear(275→128), ReLU, Dropout(0.3), Linear(128→1)

Results:

- R^2 : 0.180 ± 0.028
- MAE: 0.246 ± 0.016 km
- RMSE: 0.336 ± 0.013 km

Critical finding: Performance *degraded* compared to image-only ($\Delta R^2 = -0.099$)!

Interpretation:

- CNN features are noisy and interfere with physical features
- Naive concatenation cannot properly balance modalities
- Regression head overfits to noisy CNN features
- Demonstrates need for learned fusion mechanisms

3.2.3 Model 3: Attention Fusion

Architecture:

- Same CNN encoder (256-dim image embedding)
- Physical features: 19-dim vector
- **Cross-attention mechanism:**
 - Query: Linear(256→64) from image embedding
 - Key/Value: Linear(19→64) from physical features
 - Attention: $\alpha = \text{softmax}(QK^T / \sqrt{d_k})$
 - Output: $\text{Attn}(Q, K, V) = \alpha V$
- **Gated fusion:**

- Gate: $g = \sigma(\text{Linear}([\text{img}; \text{attn}]))$
- Fused: $f = g \odot \text{img} + (1 - g) \odot \text{attn}$
- Regressor: Linear(256→128), ReLU, Dropout(0.3), Linear(128→1)

Results:

- $R^2: 0.326 \pm 0.047$
- MAE: 0.221 ± 0.014 km (221 meters)
- RMSE: 0.304 ± 0.019 km

Analysis:

- Improves over concatenation by $\Delta R^2 = +0.146$ (81% relative improvement)
- Attention learns to downweight noisy CNN features
- Still underperforms physical baseline by $\Delta R^2 = 0.35$
- Validates hypothesis that learned fusion helps, but CNN features remain weak

3.3 Ablation Study

Comprehensive Model Comparison:

Model	R^2	MAE (km)	RMSE (km)
Physical-only GBDT	0.676	0.136	0.211
Attention Fusion CNN	0.326	0.221	0.304
Image-only CNN	0.279	0.233	0.315
Concatenation Fusion CNN	0.180	0.246	0.336

Key ablation insights:

1. **Physical vs. Image:** Physical features are $2.4\times$ stronger (R^2 gap = 0.40)
2. **Image vs. Concat:** Adding physical features *hurts* naive fusion ($\Delta R^2 = -0.10$)
3. **Concat vs. Attention:** Attention recovers 81% of lost performance ($\Delta R^2 = +0.15$)
4. **Physical vs. Best Hybrid:** Physical still wins by $2.1\times$ ($\Delta R^2 = 0.35$)

Section 4 Conclusions:

- No CNN architecture beats the physical baseline
- Root causes identified:
 - CNN trained from scratch on small dataset (933 samples)
 - No pre-training or transfer learning
 - Single-frame input (no temporal context)
 - Simple architecture (4-layer CNN insufficient)
- **Motivation for Part 5:** Pre-trained backbones + temporal modeling

4 Section 5: Advanced Deep Learning

4.1 Overview

Part 5 addresses the limitations identified in Part 4 by introducing:

1. **Pre-trained backbones:** Transfer learning from ImageNet
2. **Modern architectures:** ResNet-50, Vision Transformers (ViT)
3. **Temporal modeling:** Multi-frame sequences with attention
4. **Physics-informed losses:** Temporal consistency regularization
5. **Advanced fusion:** FiLM (Feature-wise Linear Modulation) for ERA5 integration

4.2 Part 1: Pre-Trained Backbones

4.2.1 Model: ResNet-50 Baseline

Architecture:

- **Backbone:** ResNet-50 pre-trained on ImageNet-1k (torchvision)
- **Input handling:** Grayscale images duplicated 3× to RGB channels
- **Fine-tuning strategy:**
 - Freeze: conv1, bn1, layer1, layer2, layer3 (early feature extractors)
 - Train: layer4 (high-level features) + regression head
 - Rationale: Preserve low-level edge/texture features, adapt high-level semantics
- **Regression head:** Linear(2048→512), ReLU, Dropout(0.3), Linear(512→1)

Training:

- Optimizer: AdamW (lr=1e-4, weight_decay=1e-4)
- Loss: MSE
- Batch size: 12
- Early stopping: patience=10 epochs
- Validation: Stratified 5-Fold CV

Results:

- Mean R^2 : 0.524 ± 0.046
- Mean MAE: 0.171 ± 0.012 km (171 meters)
- Mean RMSE: 0.256 ± 0.014 km
- Best fold: $R^2 = 0.621$ (Fold 0)

- Worst fold: $R^2 = 0.487$ (Fold 4)

Analysis:

- 88% improvement over Part 4 CNN (R^2 0.524 vs. 0.279)
- Pre-training on ImageNet provides strong initialization
- Still underperforms physical baseline by $\Delta R^2 = 0.15$
- ResNet's convolutional inductive bias may not be optimal for this task

4.2.2 Model: ViT-Tiny Baseline

Architecture:

- **Backbone:** ViT-Tiny (WinKawaks/vit-tiny-patch16-224) pre-trained on ImageNet-21k
- **Patch size:** 16×16 (results in 196 patches for 224×224 input)
- **Input handling:** 440×640 images resized to 224×224 , grayscale→RGB
- **Fine-tuning:** All layers trainable (ViT is small: 5.7M parameters)
- **Regression head:** Linear(192→256), ReLU, Dropout(0.3), Linear(256→1)

Training:

- Optimizer: AdamW (lr=3e-5, weight_decay=1e-4)
- Loss: MSE
- Batch size: 10
- Early stopping: patience=10 epochs
- Validation: Stratified 5-Fold CV

Results:

- Mean R^2 : 0.577 ± 0.019
- Mean MAE: 0.166 ± 0.006 km (166 meters)
- Mean RMSE: 0.241 ± 0.008 km
- Best fold: $R^2 = 0.599$ (Fold 2)
- Worst fold: $R^2 = 0.545$ (Fold 0)

Analysis:

- 107% improvement over Part 4 CNN (R^2 0.577 vs. 0.279)
- Outperforms ResNet-50 by $\Delta R^2 = +0.053$ (10% relative)
- Vision Transformer's self-attention better captures global spatial relationships

- Lower cross-fold variance ($\text{std}=0.019$) indicates more stable generalization
- Still below physical baseline by $\Delta R^2 = 0.099$

WP-1 Conclusion:

Pre-trained backbones (especially ViT) significantly improve over CNNs from scratch, but single-frame models cannot yet beat the physical baseline. As you will see in the next section, temporal modeling is what we need.

4.3 Part 2: Temporal Modeling

4.3.1 Multi-Frame Temporal ViT

Motivation:

Cloud base height evolves slowly in time (meteorological timescale: minutes to hours). A sequence of frames should provide:

- Temporal smoothing to reduce noise
- Motion/parallax cues for 3D geometry estimation
- Redundancy to handle partial occlusions

Architecture:

- **Frame encoder:** ViT-Tiny (shared weights across frames)
- **Temporal sequence:** 5 consecutive frames (center frame is target)
- **Temporal aggregation:** Multi-head self-attention (4 heads)
 - Input: Sequence of 5 frame embeddings [192-dim each]
 - Output: Aggregated temporal representation [192-dim]
- **Edge handling:** Clamp to flight boundaries (no cross-flight sequences)
- **Regression head:** Linear(192→256), ReLU, Dropout(0.3), Linear(256→1)

Training:

- Optimizer: AdamW (lr=2e-5, weight_decay=1e-4)
- Loss: MSE (on center frame only)
- Batch size: 4 (VRAM limited)
- Gradient accumulation: 4 steps (effective batch size = 16)
- Early stopping: patience=10 epochs
- Validation: Stratified 5-Fold CV

Results:

- Mean R^2 : 0.727 ± 0.052

- Mean MAE: 0.126 ± 0.008 km (126 meters)
- Mean RMSE: 0.193 ± 0.021 km
- Best fold: $R^2 = 0.823$ (Fold 4)
- Worst fold: $R^2 = 0.677$ (Fold 2)

First deep learning to beat physical baseline (!!)

- R^2 improvement: 0.727 vs. 0.668 = +8.8% relative (+0.059 absolute)
- MAE improvement: 126 m vs. 137 m = 11 meters better (8% relative)
- 26% improvement over single-frame ViT (R^2 0.727 vs. 0.577)

Analysis:

- Temporal information is critical for performance
- Multi-frame attention learns to aggregate complementary views
- Temporal smoothing reduces per-frame noise
- High variance across folds (std=0.052) suggests some domain shift remains

4.3.2 Physics-Informed Temporal Consistency Loss

Motivation:

Real cloud base height changes slowly (typically < 10 m/s vertical velocity). We can enforce this physical constraint via a temporal consistency loss:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{MSE}}(\hat{y}_{\text{center}}, y_{\text{center}}) + \lambda \cdot \mathcal{L}_{\text{temporal}} \quad (2)$$

where the temporal consistency loss penalizes rapid changes:

$$\mathcal{L}_{\text{temporal}} = \frac{1}{T-1} \sum_{t=1}^{T-1} |\hat{y}_{t+1} - \hat{y}_t| \quad (3)$$

Architecture modification:

- Model outputs predictions for all 5 frames (not just center)
- Primary loss: MSE on center frame (supervision signal)
- Regularization: Temporal consistency across all 5 predictions

Ablation study: $\lambda \in \{0.05, 0.1, 0.2\}$

Results:

Optimal configuration: $\lambda = 0.1$

Analysis:

- $\lambda = 0.05$: Under-regularized, predictions too noisy
- $\lambda = 0.10$: Optimal balance, best MAE and lowest variance

λ	Mean R^2	Mean MAE (m)	Mean RMSE (m)
0.05	0.700 ± 0.065	140.5 ± 17.0	202.2 ± 22.3
0.10	0.728 ± 0.044	125.7 ± 9.2	192.9 ± 16.8
0.20	0.728 ± 0.043	130.5 ± 12.6	192.9 ± 15.4

- $\lambda = 0.20$: Over-regularized, predictions too smooth (worse MAE)
- Temporal consistency loss provides marginal improvement: R^2 0.728 vs. 0.727
- Main benefit: Reduced cross-fold variance (std=0.044 vs. 0.052)

Best Model Summary:

Temporal ViT + Consistency Loss ($\lambda = 0.1$)

- Mean R^2 : 0.728 ± 0.044
- Mean MAE: $126 \text{ m} \pm 9 \text{ m}$
- Mean RMSE: $193 \text{ m} \pm 17 \text{ m}$

4.4 Part3: Advanced Fusion (FiLM)

4.4.1 Feature-wise Linear Modulation

Motivation:

Part 4's attention fusion ($R^2 = 0.326$) failed to effectively integrate ERA5 features. FiLM (Feature-wise Linear Modulation) is a more expressive fusion mechanism:

$$\text{FiLM}(x, z) = \gamma(z) \odot x + \beta(z) \quad (4)$$

where x is the image features, z is the ERA5 features, and γ, β are learned parameters.

Architecture:

- **Image encoder:** ViT-Tiny (192-dim embeddings)
- **ERA5 encoder:** MLP: Linear(9→64), ReLU, LayerNorm, Linear(64→64)
- **FiLM generator:**
 - Input: 64-dim ERA5 encoding
 - Gamma branch: Linear(64→192), Sigmoid (scale near 1.0)
 - Beta branch: Linear(64→192), no activation
- **Modulation:** $x_{\text{fused}} = \gamma \odot x_{\text{image}} + \beta$
- **Regression head:** Linear(192→256), ReLU, Dropout(0.3), Linear(256→1)

Implementation details:

- ERA5 features normalized (z-score) before encoding
- Gamma gating (sigmoid) ensures $\gamma \approx 1.0$ initially (identity initialization)

- Gradient clipping (max_norm=1.0) for training stability
- LayerNorm in FiLM generator to prevent exploding activations

Training:

- Optimizer: AdamW (lr=3e-5, weight_decay=1e-4)
- Loss: MSE
- Batch size: 10
- Early stopping: patience=10 epochs
- Validation: Stratified 5-Fold CV

Results:

- Mean R^2 : 0.542 ± 0.026
- Mean MAE: 0.166 ± 0.006 km (166 meters)
- Mean RMSE: 0.251 ± 0.009 km
- Best fold: $R^2 = 0.589$ (Fold 1)
- Worst fold: $R^2 = 0.507$ (Fold 0)

Analysis:

- FiLM underperforms image-only ViT (R^2 0.542 vs. 0.577, $\Delta = -0.035$)
- ERA5 features do not improve single-frame models
- Possible explanations:
 - ERA5 spatial resolution (25 km) too coarse for imagery (200 m pixels)
 - Hourly temporal resolution misses sub-hourly cloud dynamics
 - Atmospheric state may not strongly constrain cloud base height
 - Better fusion may require cross-modal attention (not just affine modulation)
- Training was stable (no gradient explosions after fixes)

WP-3 Conclusion:

FiLM fusion validates proper integration of ERA5 features but does not improve performance. Atmospheric features may be most useful in physical models (GBDT), not deep learning. Future work should explore cross-modal attention.

Rank	Model	R^2	MAE (m)	RMSE (m)
1	Temporal ViT + Consistency	0.728	126	193
2	Temporal ViT	0.727	126	193
-	Physical GBDT (baseline)	0.668	137	213
3	ViT-Tiny	0.577	166	241
4	FiLM Fusion	0.542	166	251
5	ResNet-50	0.524	171	256

4.5 Part 5 Summary

Model Performance Ranking:

Key findings:

1. **Temporal modeling is critical:** +26% R^2 improvement over single-frame
2. **ViT outperforms ResNet:** Attention-based architectures superior for this task
3. **Physics-informed losses help marginally:** Temporal consistency reduces variance
4. **ERA5 fusion remains unsolved:** FiLM does not improve over image-only
5. **Production recommendation:** Temporal ViT + Consistency ($\lambda = 0.1$)

5 Cross-Part/Task Analysis

5.1 Performance Evolution

R^2 progression across test:

- Part 3: Physical baseline = 0.668
- Part 4: Best CNN (attention) = 0.326 (-51% vs. baseline)
- Part 5: Pre-trained ViT = 0.577 (-14% vs. baseline)
- Part 5: Temporal ViT = 0.727 (+9% vs. baseline) ✓

Key insights:

1. **Transfer learning is essential:** Pre-training on ImageNet provides critical initialization for small datasets (933 samples)
2. **Temporal context is the breakthrough:** Single-frame models cannot compete with multi-frame temporal attention
3. **Architecture matters:** ViT (self-attention) outperforms ResNet (convolution) by 10% for this task
4. **Physical features excel in classical ML:** GBDT effectively combines geometric + atmospheric features, but deep learning struggles to fuse modalities
5. **Dataset size is a bottleneck:** 933 samples is small for deep learning; pre-training mitigates but does not eliminate this limitation

5.2 Feature Importance: Physical vs. Learned

Physical GBDT top features:

1. BLH (boundary layer height)
2. Shadow-derived geometric height
3. Solar zenith angle
4. LCL (lifting condensation level)

Deep learning attention maps (qualitative):

- ViT looks to cloud edges and shadow boundaries
- Temporal ViT shows strong attention to motion/parallax
- Limited attention to atmospheric context (consistent with FiLM failure)

Conclusion: Physical models excel at explicit feature engineering, while deep learning discovers geometric features implicitly (but requires temporal context to succeed).

Model	Parameters	Train Time/Fold	Inference (ms)	VRAM (GB)
Physical GBDT	–	2 min	0.1	–
CNN (Section 4)	1.2M	15 min	5	2
ResNet-50	23.5M	45 min	15	4
ViT-Tiny	5.7M	30 min	8	3
Temporal ViT	6.1M	60 min	40	6

5.3 Computational Requirements

Production considerations:

- **Physical GBDT:** Fastest, lowest resource, production-ready
- **Temporal ViT:** Best accuracy, but $400\times$ slower inference
- **Trade-off:** $+8.8\% R^2$ improvement costs $400\times$ latency
- **Recommendation:** Deploy both models (GBDT for real-time, ViT for post-processing)

5.4 Validation Protocol Robustness

Stratified K-Fold CV stability:

All models trained with stratified 5-fold CV (CBH target binned into 10 quantiles). Cross-fold variance analysis:

Model	Mean R^2	Std R^2
Physical GBDT	0.668	0.044
ViT-Tiny	0.577	0.019
Temporal ViT	0.727	0.052
Temporal + Consistency	0.728	0.044

Observation: Temporal consistency loss reduces variance ($0.052 \rightarrow 0.044$), matching GBDT stability.

Domain shift (Flight F4):

Flight F4 has anomalously low CBH (mean = 0.697 km vs. 0.917 km for F0). Leave-One-Flight-Out CV on F4 yields catastrophic failure ($R^2 = -3.13$), confirming extreme domain shift. Stratified K-Fold CV is the correct validation protocol for model development.

6 Data

6.1 ERA5 Processing Details

File: sow_outputs/wp2_atmospheric/WP2_Features.hdf5

Processing pipeline:

1. ERA5 NetCDF files downloaded from ECMWF (119 daily files, Oct 2024 – Feb 2025)
2. Spatiotemporal matching: Nearest neighbor (lat/lon) + nearest hourly timestamp
3. Feature extraction: BLH, temperature, dewpoint, pressure (surface + levels)
4. Derived features: LCL, inversion height, stability index, moisture gradient
5. Quality control: 933/933 samples successfully matched (100%)

Verification:

- Mean BLH = 658 m (physically realistic for marine boundary layer)
- LCL correlates with observed CBH ($r = 0.42$)
- Stability index consistent with stratocumulus regime (3.81 K/km)
- All values within expected ranges for coastal/oceanic conditions

7 Deployment Recommendations

7.1 Recommended Model Architecture

Primary Model: Temporal ViT with Temporal Consistency Loss ($\lambda = 0.1$)

Performance:

- $R^2 = 0.728$ (vs. 0.668 baseline = +9% improvement)
- MAE = 126 m (vs. 137 m baseline = 11 m better)
- RMSE = 193 m (vs. 213 m baseline = 20 m better)

Architecture details:

- Frame encoder: ViT-Tiny (WinKawaks/vit-tiny-patch16-224)
- Input: 5-frame sequence (temporal window ≈ 10 seconds at 2 Hz)
- Temporal aggregation: 4-head self-attention
- Loss: MSE + $0.1 \times$ temporal consistency
- Parameters: 6.1M (deployable on edge devices)

Model checkpoints:

- 5 fold models saved: `sow_outputs/wp5/models/temporal_consistency/lambda_0.1/fold_*.pth`
- Best fold: Fold 1 ($R^2 = 0.782$)
- Recommended for deployment: Ensemble of all 5 folds (vote averaging)

7.2 Alternative: Dual-Model Deployment

Operational strategy:

Deploy both Physical GBDT and Temporal ViT in production:

- **Real-time processing:** Physical GBDT (0.1 ms inference, CPU-only)
 - Use for immediate CBH estimates during flight
 - Requires only geometric features (shadow detection)
 - Fallback when temporal context unavailable (e.g., first frames)
- **Post-processing:** Temporal ViT (40 ms inference, GPU recommended)
 - Apply after flight for highest-accuracy retrievals
 - Requires 5-frame buffer (2.5 seconds at 2 Hz)
 - Primary output for scientific analysis

7.3 Known Limitations and Mitigations

Limitation 1: Domain shift (Flight F4)

- **Issue:** F4 has anomalously low CBH (mean = 0.697 km vs. 0.917 km)
- **Impact:** LOO CV on F4 yields $R^2 = -3.13$ (catastrophic failure)
- **Mitigation:** Stratified K-Fold CV for development; few-shot fine-tuning for deployment on new flight regimes

Limitation 2: Temporal edge cases

- **Issue:** First 2 frames of each flight lack full 5-frame context
- **Mitigation:** Fallback to Physical GBDT for edge frames; use padding/extrapolation

Limitation 3: Computational cost

- **Issue:** Temporal ViT is 400× slower than GBDT (40 ms vs. 0.1 ms)
- **Mitigation:** Deploy dual-model system (GBDT real-time, ViT post-processing)

Limitation 4: ERA5 fusion failure

- **Issue:** FiLM fusion ($R^2 = 0.542$) underperforms image-only ViT ($R^2 = 0.577$)
- **Mitigation:** Use ERA5 features only in GBDT; explore cross-modal attention in future work

8 Future Work

8.1 High Priority

1. Uncertainty Quantification

- Implement Monte Carlo Dropout for prediction confidence estimates
- Calibrate uncertainty to some operational requirements (e.g., 90% confidence intervals)
- Use uncertainty to flag low-confidence predictions for manual review (can tie into current paper)

2. Cross-Modal Attention

- Implement Query(image) \times Key/Value(ERA5) attention
- Compare to FiLM fusion (current best $R^2 = 0.542$)
- Hypothesis: Explicit attention to atmospheric features may improve over just modulation

3. Domain Adaptation for Flight F4

- Few-shot fine-tuning on low-CBH regimes
- Meta-learning for rapid adaptation to new flight conditions
- Investigate domain-invariant representations

4. Ensemble Methods

- Combine Temporal ViT + Physical GBDT (stacking or weighted averaging)
- Expected performance: $R^2 > 0.74$
- Improved robustness and uncertainty estimates

8.2 Medium Priority (Research Extensions)

5. Model Compression and Optimization

- TorchScript / ONNX export for deployment
- Quantization (INT8) for faster inference
- Knowledge distillation: Temporal ViT \rightarrow some smaller student model
- Target: 10 \times speedup with < 5% accuracy loss (only really good to do if you want real time inference for during flight)

6. Explainability and Visualization

- Attention map visualization (which image regions drive predictions?)
- Temporal attention flow (how does model integrate multi-frame information?)

7. Self-Supervised Pre-Training

- Revisit earlier research on self-supervised learning (MAE, contrastive learning)
- Pre-train ViT on unlabeled ER-2 imagery (10,000+ frames available)
- Fine-tune on labeled CBH data (933 samples)
- Hypothesis: Domain-specific pre-training may improve over ImageNet initialization

8.3 Low Priority (Exploratory Research)

8. Alternative Architectures

- Mamba / S4 (state-space models for efficient sequence modeling)
- Video Swin Transformer (hierarchical spatiotemporal attention)
- 3D CNNs (spatiotemporal convolution)

9. Multi-Task Learning

- Joint prediction: CBH + cloud optical depth + cloud top height
- Auxiliary tasks: Shadow detection, cloud segmentation
- Hypothesis: Shared representations may improve generalization

10. Just Get More Data

- Just test on the new acquisition system Drew made - ideally would give better predictions but also very different data from what we have now

9 Conclusions

9.1 Summary of Achievements

1. **Physical baseline established:** XGBoost GBDT with geometric and atmospheric features achieves $R^2 = 0.668$, MAE = 137 m
2. **Hybrid CNN development:** Attention fusion outperforms naive concatenation but cannot beat physical baseline
3. **Pre-trained backbones:** ResNet-50 and ViT-Tiny significantly improve over CNNs from scratch
4. **Breakthrough with temporal modeling:** Temporal ViT is the first deep learning model to beat the physical baseline, achieving $R^2 = 0.727$, MAE = 126 m
5. **Physics-informed regularization:** Temporal consistency loss (Task 2.2) provides marginal improvement and reduced variance
6. **Production-ready model:** Temporal ViT + Consistency Loss ($\lambda = 0.1$) recommended for deployment with $R^2 = 0.728$, MAE = 126 m

9.2 Critical Findings

- **Temporal information is essential:** Multi-frame models outperform single-frame by 26% R^2
- **Vision Transformers excel:** ViT systematically outperforms ResNet for this task
- **Transfer learning critical:** ImageNet pre-training provides crucial initialization for small datasets
- **ERA5 fusion remains challenging:** Atmospheric features improve GBDT but not deep learning models

9.3 Final Remarks

The transition from physical baselines ($R^2 = 0.668$) to advanced deep learning with temporal modeling ($R^2 = 0.728$) represents a **9% performance improvement** and validates the deep learning approach for cloud base height retrieval.

The Temporal ViT model is recommended for production deployment pending offline validation and pilot testing. Future work should focus on uncertainty quantification, domain adaptation, and ensemble methods to further improve robustness and operational readiness. (real test will be when I feed in new data to this model)