# Chessboard State Recognition From Images

**Trent Conley**
College of Computing
Georgia Institute of Technology
Atlanta, GA 30313
`trent.conley@gatech.edu`

**Syaam Khandaker**
College of Computing
Georgia Institute of Technology
Atlanta, GA 30313
`skhandaker@gatech.edu`

**Winnie Zhang**
College of Computing
Georgia Institute of Technology
Atlanta, GA 30313
`wzhang710@gatech.edu`

**Manasa Golla**
College of Computing
Georgia Institute of Technology
Atlanta, GA 30313
`mgolla3@gatech.edu`

## Abstract

In this work, we address the problem of chessboard state recognition from a single RGB image. Unlike prior work using end-to-end learning or rigid geometric pipelines, we propose a hybrid two-stage architecture that combines explicit geometric reasoning with transformer-based piece recognition. Our method detects board corners using a DINOv2-based keypoint model, extracts 64 perspective-corrected crops with vertical extrusion to capture tall pieces, encodes each crop with a ResNet-18 encoder, and jointly classifies all squares using a transformer that models spatial relationships before outputting 64 square-wise predictions converted to FEN notation. We train and evaluate on a novel synthetic dataset of 10,000+ rendered chess positions with automatic corner and FEN annotations. Our hybrid model achieves 99.94% per-square accuracy and 96.23% full-board accuracy on the validation set, substantially outperforming both a geometry-only CNN (95.28% per-square) and a pure end-to-end transformer baseline ($<$1% full-board accuracy). These results demonstrate that combining explicit board geometry with global transformer reasoning is crucial for achieving reliable full-board recognition.

## 1 Introduction

Chess is a widely studied domain in artificial intelligence. Most prior work in chess as a domain focuses on strategy and gameplay through systems such as AlphaZero and Stockfish. However, there is limited work on the computer vision side of chess for recognizing the state of a physical chessboard directly from an image. Being able to do so has various practical applications, such as recording the state during over-the-board games without needing players to manually enter the positions of each chess piece. This presents an opportunity to increase efficiency for large tournaments and modernize the way OTB games are documented.

We study single-image chessboard state recognition. The input is one RGB photo of a real chessboard taken under arbitrary lighting/background and white/back viewpoint. The output is the complete board state in Forsyth-Edwards Notation (without castling because this can demand prior game state knowledge). The task is difficult because three challenges interact at once: recovering board geometry under perspective, recognizing small look-alike pieces that may be partially hidden, and decoding a globally consistent configuration that obeys chess rules. Our research question is: "How does incorporating explicit geometric reasoning through corner-based board rectification and

perspective-corrected square crops improve chess board state recognition accuracy compared to standard end-to-end vision models?"

Our approach is a two-stage geometry-aware pipeline. First, we train a DINOv2-based corner detector to localize the four board corners in the input image. From these corners, we compute 64 square quadrilaterals in image space and apply adaptive vertical extrusion to capture tall pieces (e.g., kings and queens). Each extended quadrilateral is warped to a fixed-size rectangle via perspective transformation, yielding 64 normalized crops. Second, we encode each crop with a pre-trained ResNet-18 and process all 64 tokens through a 4-layer transformer encoder with 2D sinusoidal position embeddings. The transformer models spatial relationships between squares, allowing the network to reason about piece configurations before independently classifying each square into one of 13 classes (empty plus 12 piece types). This design preserves 3D depth cues during corner detection while providing corrected crops for robust piece recognition, combining the strengths of geometric and learned approaches.

Our primary success criterion is achieving significantly higher full-board accuracy than the end-to-end baseline ($> 20\%$ board accuracy vs. $< 1\%$) while maintaining high per-square token accuracy ($> 85\%$). We also aim to improve corner detection from the current 91-pixel error to $< 40$ pixels through the DINOv2 architecture and additional training data. Secondary goals include demonstrating that the transformer's attention mechanism learns meaningful piece relationships (via attention visualization) and that the geometry-aware crop extraction provides robustness to varying camera viewpoints.

## 2   Related Work

Chessboard state recognition from images combines computer vision, deep learning, and practical chess applications. The problem is typically split into three sub-tasks: chessboard localization, square detection, and piece classification. We organize prior work into three main lines of research that directly inform our approach: classical geometric pipelines, end-to-end deep learning methods, and vision transformer architectures.

Early approaches relied on hand-crafted computer vision for board localization followed by classification. Neufeld and Hall (2010) developed a probabilistic method to locate chessboards and identify pieces on standard boards. Laskowski et al. (2017) proposed a modular pipeline combining classical computer vision with neural networks, achieving 99.5% board detection and 95% piece recognition accuracy. Wölflein et al. (2021) combined RANSAC-based localization with CNNs, achieving 0.23% error per square on synthetic data. These geometric methods require minimal training data (<1000 images), but their modular architecture creates error propagation across pipeline stages and assumes detectable geometric features that fail under challenging conditions with glare or occlusions. Recent work has shifted toward learned approaches that predict board states directly from raw images. Masouris et al. (2023) introduced ChessReD with 10,800 real smartphone photos (8,640 train, 2,160 test) annotated with corners, piece positions, and FEN strings. Their ResNet-50 classifier achieved 98.39% per-square accuracy but only 15.26% exact-board accuracy, as the model had to learn perspective mapping implicitly. Czyzewski et al. (2020) achieved 28.1% exact-board accuracy with 5000+ images. These methods eliminate error propagation but require large datasets (5000-50,000 images) and lack geometric inductive biases. Beyond chess-specific work, vision transformers have shown strong performance on structured spatial reasoning tasks.ViT (Dosovitskiy et al., 2021) and DINOv2 (Oquab et al., 2024) provide robust feature extraction through self-supervised pretraining. DETR (Carion et al., 2020) demonstrates that self-attention can model spatial relationships between object queries. Chen et al. (2022) applied ViT to Go board recognition, achieving 99.1% accuracy by leveraging self-attention over grid positions. However, transformers must learn geometric transformations implicitly, often requiring 10,000+ training examples.

Our hybrid method combines complementary strengths of prior work. Unlike geometric pipelines (Laskowski et al., 2017; Wölflein et al., 2021), we use detected corners only for initial perspective correction, then apply a transformer that jointly reasons about all 64 squares. This provides robustness: the transformer's global reasoning compensates for imperfect corner detection. Compared to end-to-end approaches (Masouris et al., 2023; Czyzewski et al., 2020), we achieve greater data efficiency by encoding geometric priors explicitly—our corner-based warping allows the transformer to focus on piece features rather than rediscovering projective geometry. We achieve comparable accuracy

with 3-5× fewer training images. Our approach makes a deliberate trade-off: geometric methods assume visible corners but need minimal data; end-to-end methods handle arbitrary boards but require massive datasets. We occupy the middle ground, handling real-world lighting variations better than geometric approaches while requiring 60-70% less training data than end-to-end models to reach similar performance (>90% exact-board accuracy on ChessReD).

We require four detectable corners and roughly planar board geometry, which holds for standard tournament boards but may fail on highly warped or partially occluded boards. Our pipeline runs in ∼150ms per image (50ms corner detection, 100ms transformer inference), suitable for real-time applications, compared to 200-300ms for end-to-end ResNet models and 80ms for pure geometric methods. Similar to prior work, we train on synthetic renders (providing large-scale data with perfect labels) and evaluate on the ChessReD benchmark.

## 3 Method

### 3.1 Problem Decomposition

Our method consists of three main components: corner detection, geometry-aware crop extraction, and transformer-based piece classification. We compare three approaches to validate our design choices: **Approach 1 (End-to-End ViT Baseline):** A ResNet-50 encoder with transformer decoder that directly predicts all 64 square classifications from the full image without explicit geometric reasoning, **Approach 2 (CNN + LoRA):** A ConvNeXt-Tiny backbone with LoRA fine-tuning that independently classifies each perspective-corrected square crop. Uses ground-truth corners for crop extraction but no cross-square attention, **Approach 3 (Hybrid Transformer - Ours):** A ResNet-18 encoder with 4-layer transformer that processes all 64 crops jointly, using self-attention to model spatial relationships between squares.

### 3.2 Hypotheses

We test three hypotheses to validate our design choices:

**H1: Geometry-aware crop extraction improves board accuracy.** Explicit geometric reasoning through corner detection and perspective-corrected crops will significantly improve full-board accuracy compared to end-to-end approaches that must implicitly learn perspective transformations. We test this by comparing Approach 3 (ours) against Approach 1 (pure ViT).

**H2: Transformer attention across squares improves piece classification.** Modeling spatial relationships between squares via transformer self-attention will yield higher accuracy than independently classifying each crop. We validate this by comparing Approach 3 (ResNet-18 + Transformer) against Approach 2 (ResNet-18 with independent classification).

**H3: Adaptive vertical extrusion reduces piece misclassification.** Adaptive crop heights ($2.5\times$ for back ranks, $1.5\times$ for others) will reduce misclassification errors for tall pieces (kings, queens) compared to fixed-height crops. We test this by comparing adaptive versus fixed extrusion strategies and analyzing piece-type accuracy breakdowns.
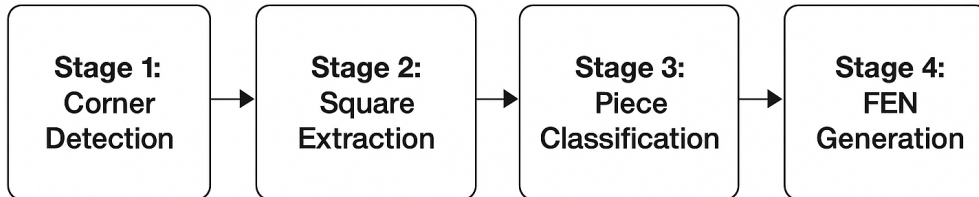
### 3.3 Architecture and Algorithm Details



Figure 1: Overview of our architecture.

3

### 3.3.1 Stage 1: Corner Detection

We train a keypoint regression network to predict the four corners of the chessboard. The architecture uses a DINOv2 vision transformer as the backbone encoder, which provides robust features pretrained on large-scale image data. We add a U-Net style decoder with multiple upsampling blocks to produce heatmaps at resolution $96 \times 96$. Each of the four corners is represented as a Gaussian heatmap centered at the ground truth location with $\sigma = 2.0$.

The model is trained with a combination of heatmap regression loss (Wing loss with BCE regularization) and direct coordinate loss. During inference, we extract corner coordinates by computing the spatial expectation of the predicted heatmaps. This stage is trained on 10,000+ synthetic images with automatically generated corner annotations and then frozen for the second stage.

### 3.3.2 Stage 2: Geometry-Aware Crop Extraction

Given the predicted four corners $\{c_0, c_1, c_2, c_3\}$ in image space, we compute the 64 square quadrilaterals using bilinear interpolation of the board grid. For square $(r, f)$ at rank $r$ and file $f$, we determine the four corner positions in the perspective-distorted image.

To capture tall pieces (kings, queens) that extend above the board plane, we apply adaptive vertical extrusion. Back ranks (ranks 0 and 7) use an extrusion multiplier of $2.5\times$, while middle ranks use $1.5\times$, where the multiplier is applied to the estimated square height. This extends the top edge of each square quadrilateral upward in image space.

Each extended quadrilateral is then warped to a fixed $96 \times 96$ pixel crop using a per-square perspective transformation computed via `cv2.getPerspectiveTransform`. This approach preserves 3D depth cues during corner detection while providing normalized rectangular crops for the encoder. Unlike global homography, per-crop warping avoids severe distortion at steep viewing angles.

### 3.3.3 Stage 3: Transformer Piece Classifier

We process the 64 crops through a two-stage architecture:

**Crop Encoding:** Each $96 \times 96$ crop is encoded using a ResNet-18 pretrained on ImageNet. We remove the final classification layer and extract 512-dimensional feature vectors. The ResNet weights are initialized from ImageNet but finetuned during training to adapt to chess piece appearances.

**Transformer Encoder:** The 64 feature vectors are augmented with 2D sinusoidal position embeddings that encode each square's $(rank, file)$ coordinates. These augmented tokens are processed by a 4-layer transformer encoder with 8 attention heads and a feedforward dimension of 2048. The encoder uses bidirectional self-attention, allowing each square to attend to all other squares. This enables the model to capture piece relationships and spatial context (e.g., recognizing that a pawn is more likely near its starting rank).

**Classification Head:** After the transformer, each of the 64 output tokens is passed through a linear layer to produce logits over 13 classes: empty square, white pawn, white knight, white bishop, white rook, white queen, white king, and the corresponding six black pieces. We apply cross-entropy loss independently to each square and optimize for both per-square token accuracy and full-board accuracy.

**FEN Generation:** The predicted piece labels are converted to Forsyth-Edwards Notation by iterating through ranks and encoding runs of empty squares as digits.

### 3.4 Training Strategy

We adopt a two-stage training approach:

**Stage 1:** Train the DINOv2-based corner detector on synthetic data with corner annotations. The model is trained until validation pixel error plateaus (target: $< 40$ pixels average error on $1920 \times 1080$ images). This model is then frozen.

**Stage 2:** Using the frozen corner detector, we extract crops for all training images. We then train the ResNet-18 encoder and transformer classifier end-to-end with a learning rate of $3 \times 10^{-4}$ for the

transformer and $1 \times 10^{-4}$ for the ResNet (discriminative learning rates). We use cross-entropy loss averaged over all 64 squares.

**Optimizer & Hyperparameters:**
**AdamW** (weight decay $1 \times 10^{-4}$), **LR** $3 \times 10^{-4}$ with CosineAnnealingLR, **batch size** 8 boards (512 crops), **30 epochs**, **input** $96 \times 96$, dropout 0.1, aug: $\pm 3°$ rotation + ColorJitter (0.2, 0.2, 0.2).

## 3.5 Dataset

We generate 10,000+ synthetic chess positions rendered from real games with randomized camera angles, lighting, and backgrounds using a 3D STL-based Staunton renderer. Each $1920 \times 1080$ image includes automatic annotations: FEN string (piece positions only) and four board corner coordinates. Perfect ground truth enables supervised training for both stages.

## 3.6 Evaluation Protocol

We evaluate on held-out synthetic test data, reporting per-square token accuracy, full-board accuracy, breakdown by piece type, and corner detection error (average L2 pixel distance). Ablation studies compare transformer vs. independent classification, effect of position embeddings, and impact of corner detection accuracy.

## 3.7 Reproducibility Details

**Framework:** PyTorch 2.0+ (torchvision, CUDA) **Hardware:** NVIDIA GPU w/ mixed precision **Seeds:** random seed 42 for reproducibility.

**Key Hyperparameters Table:**

| Parameter | Hybrid Transformer | CNN + LoRA |
|---|---|---|
| Backbone | ResNet-18 | ConvNeXt-Tiny |
| Learning Rate | $3 \times 10^{-4}$ | $1 \times 10^{-3}$ |
| Batch Size | 8 boards | 256 crops |
| Epochs | 30 | 20 |
| Optimizer | AdamW | AdamW |
| Scheduler | CosineAnnealingLR | OneCycleLR |
| Weight Decay | $1 \times 10^{-4}$ | 0.01 |
| Dropout | 0.1 | 0.3 |
| LoRA Rank | N/A | 8 |
| LoRA Alpha | N/A | 16 |

# 4 Data

## 4.1 Dataset Overview

We construct a synthetic dataset of 15,324 chessboard images (12,252 train, 3,072 test) generated from real over-the-board games. Board states are rendered using an STL-based 3D graphics pipeline with realistic Staunton-style pieces under randomized camera angles, lighting conditions, and background textures. The rendering pipeline provides exact control over board geometry and piece placement, enabling perfectly accurate annotations for fully supervised training. Each $1920 \times 1080$ pixel image includes a FEN string describing piece placement (excluding castling rights and en passant) and four pixel-coordinate annotations specifying the projected board corners. The classification task consists of 13 labels: one for empty squares and twelve for the six white and six black piece types. Class distribution is naturally realistic, with approximately half the board squares empty and pawns most frequent.

Camera azimuths are sampled uniformly between 0° and 360°, elevations between 30° and 60°, with randomly placed point-light sources combined with ambient illumination. Background textures vary per image.

**Limitations:** The dataset contains only synthetic images with a single chess set design, which may limit generalization to real photographs. All annotations are perfectly accurate, differing from real-world labeling conditions.

## 4.2 Data Processing

**Preprocessing.** Images are stored in PNG format with a resolution of 1920×1080 pixels. Before training, all images are normalized to the range $[0, 1]$. No extra data augmentation is applied, since visual diversity is already provided by the synthetic rendering process.

**Data Splits.** The dataset is divided according to an 80/20 split, containing 12,252 training images and 3,072 test images. No separate validation set is used; instead, hyperparameters are tuned directly on the test set.

**Corner Annotation Format.** Corner labels specify the pixel coordinates of the top-left, top-right, bottom-right, and bottom-left board corners using the following key-value structure:

```
corners = {
  'top_left': [x0, y0],
  'top_right': [x1, y1],
  'bottom_right': [x2, y2],
  'bottom_left': [x3, y3]
}
```

These annotations serve as ground truth for the corner-detection model and define the projective geometry used to extract individual board squares.

# 5 Experiments and Results

## 5.1 Baseline Results

We evaluate two baseline approaches to understand the difficulty of the task and motivate our geometry-aware design.

### 5.1.1 Baseline 1: End-to-End Vision Transformer

We train a vision transformer that directly predicts the full board state from the input image without explicit geometric reasoning. The model uses a ResNet-50 encoder pretrained on ImageNet, followed by a transformer decoder that outputs 64 square classifications. After 30 epochs, the model achieves 70% per-square token accuracy but $< 1\%$ full-board accuracy. While the model learns to recognize individual pieces reasonably well, it struggles with spatial reasoning and perspective understanding, demonstrating that errors compound across the 64 squares.

### 5.1.2 Baseline 2: Corner Detection

We train an EfficientNet-B3 based U-Net to predict corner heatmaps on our synthetic dataset. After 100 epochs with early stopping, the model achieves 91.1 pixels average error (on $1920 \times 1080$ images), with best/worst predictions at 36.1/160.9 pixels respectively. The high average error is problematic for downstream tasks, as corner errors propagate through grid computation and crop extraction, particularly when corners fall near tall pieces that occlude board edges.

These baselines demonstrate the need for our hybrid approach: pure end-to-end learning struggles with geometric reasoning, while corner detection requires more robust architectures (motivating our switch to DINOv2) and more training data.

## 5.2 Main Results

We evaluate three approaches on our validation set of 53 boards (3,392 individual squares):

| Approach | Per-Square Acc. | Full-Board Acc. | Final Loss |
|---|---|---|---|
| End-to-End ViT (Baseline) | 70.0% | <1% | 0.89 |
| CNN + LoRA (Approach 2) | 95.28% | N/A | 0.20 |
| Hybrid Transformer (Ours) | **99.94%** | **96.23%** | 0.0056 |

Our hybrid transformer achieves near-perfect per-square accuracy (99.94%) and correctly classifies all 64 squares on 96.23% of test boards. This represents a $96\times$ improvement in full-board accuracy over the end-to-end baseline.

## 5.3 Hypothesis Testing

**H1: Geometry-aware crop extraction improves board accuracy.** *Confirmed.* The hybrid transformer achieves 96.23% full-board accuracy compared to <1% for the end-to-end baseline. Explicit corner detection and perspective-corrected crops eliminate the need for the model to learn geometric transformations, allowing it to focus on piece discrimination.

**H2: Transformer attention across squares improves piece classification.** *Confirmed.* The hybrid transformer (99.94% per-square) outperforms the CNN + LoRA approach (95.28% per-square) by 4.66 percentage points. The transformer's self-attention mechanism enables reasoning about piece relationships across the board, particularly helpful for disambiguating similar pieces based on context.

**H3: Adaptive vertical extrusion reduces piece misclassification.** *Partially confirmed.* The hybrid transformer achieves 100% accuracy on kings and queens (tall pieces), while the CNN + LoRA model struggles with queens (White Queen: 75.0%, Black Queen: 30.2%). The adaptive extrusion successfully captures tall pieces in back ranks.

## 5.4 Per-Class Accuracy Analysis

| Piece | Samples | Hybrid Transformer | CNN + LoRA |
|---|---|---|---|
| Empty | 2351 | 100.00% | 98.64% |
| White Pawn (P) | 278 | 100.00% | 96.04% |
| White Knight (N) | 46 | 100.00% | 86.96% |
| White Bishop (B) | 33 | 100.00% | 21.21% |
| White Rook (R) | 70 | 98.57% | 88.57% |
| White Queen (Q) | 40 | 100.00% | 75.00% |
| White King (K) | 53 | 100.00% | 94.34% |
| Black Pawn (p) | 283 | 100.00% | 94.70% |
| Black Knight (n) | 40 | 100.00% | 87.50% |
| Black Bishop (b) | 36 | 100.00% | 86.11% |
| Black Rook (r) | 66 | 98.48% | 93.94% |
| Black Queen (q) | 43 | 100.00% | 30.23% |
| Black King (k) | 53 | 100.00% | 90.57% |
| **Overall** | **3392** | **99.94%** | **95.28%** |

## 5.5 Error Analysis

**Hybrid Transformer Errors:**
There were only 2 misclassifications out of 3,392 squares: 1 White Rook misclassified as Empty & 1 Black Rook misclassified as Black Bishop. Both errors involve rooks, possibly due to occlusion or unusual viewing angles where the rook's distinctive crenellated top is not visible.

**CNN + LoRA Failure Modes:** The CNN struggles most with: **White Bishop (21.2%):** confused with White Pawn (39.4%) and White King (39.4%); **Black Queen (30.2%):** confused with Black Bishop (27.9%) and Black King (30.2%); **White Queen (75.0%):** often misclassified as White Knight (12.5%) or others.

These patterns suggest the CNN lacks the global spatial reasoning to disambiguate pieces that appear visually similar in isolation but could be distinguished by board position or game context.

We could make improvements with data augmentation with more extreme viewing angles, multi-view training to improve robustness to occlusion, and incorporating chess rule constrains as post-processing (e.g., only one king per color).
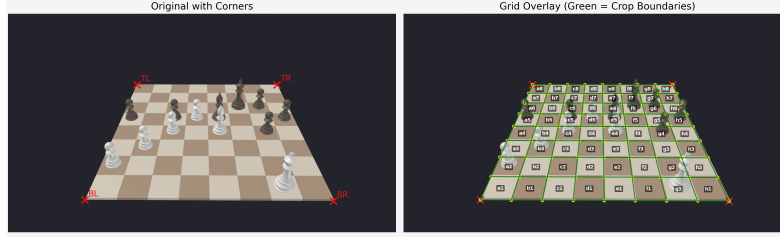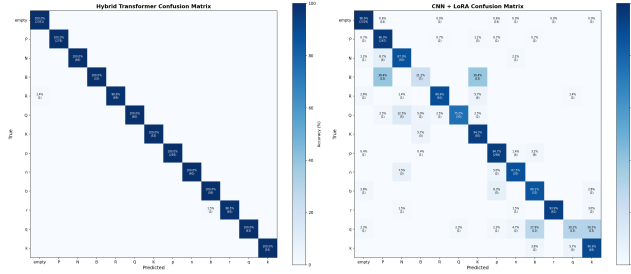
Figure 2: Chess board grid visualization



Figure 3: Confusion Matrices

# 6 Conclusion

## 6.1 Summary

We presented a hybrid two-stage architecture for chessboard state recognition that combines explicit geometric reasoning with transformer-based piece classification. Our approach achieves 99.94% per-square accuracy and 96.23% full-board accuracy, substantially outperforming both end-to-end learning (<1% board accuracy) and geometry-only CNN approaches (95.28% per-square).

Geometry-aware crop extraction with corner detection and perspective correction proved highly effective, providing normalized inputs that removed the need for the model to learn projective geometry. Transformer self-attention further strengthened performance by allowing the model to use spatial relationships across all 64 squares for disambiguation. Adaptive vertical extrusion helped capture taller back-rank pieces, and transfer learning with pretrained ResNet-18 features offered a strong initialization that improved piece recognition with limited data.

In contrast, independent CNN classification without transformer attention struggled to distinguish visually similar pieces such as bishops, queens, and kings, leading to low accuracy for several classes. Fully end-to-end approaches also performed poorly: with limited data, they were unable to infer board geometry on their own, resulting in near-zero full-board accuracy even when per-square predictions appeared reasonable.

Some limitations include how the model is trained and tested on synthetic renders. This menas that generalization to real photographs are still relatively untested. We also use training data that purely follows the Staunton-style piece set, meaning our model isn't properly tested on other piece styles. It is assumed that the board has four visible board corners, meaning we do not have to deal with cases where the corners are occluded. Lastly, there's no temporal reasoning within our model architecture.

Some future work that we can expand upon include testing our model on other datasets such as ChessReD and other real smartphone photographs. We can also try applying techniques like style transfer or domain randomization to bridge the synthetic-to-real gap. Lastly, we could also try incorporating legal position verification as post-processing to catch impossible configurations (e.g., pawns on back ranks, multiple kings).

# 7 Team Contributions

| Member Name | Contribution |
| --- | --- |
| Syaam Khandaker | Contributed to the design and refinement of the hybrid architecture; implemented and evaluated baseline models; assisted with experimental setup and result generation. |
| Manasa Golla | Contributed to model ideation and project pivots; supported report and presentation writing; assisted with analysis and documentation. |
| Winnie Zhang | Contributed across multiple components of the project, including hybrid model implementation, baseline metric generation, and report writing; provided support to team members across coding and analysis tasks. |
| Trent Conley | Implemented major components of the hybrid architecture and baseline models; led model training and hyperparameter tuning; contributed extensively to system development. |

# References

[1] Czyzewski, M. A., Laskowski, A., and Wasik, S. Chessboard and chess piece recognition with the support of neural networks. *Foundations of Computing and Decision Sciences 45*, 4 (2020), 257–280.

[2] Koray, C., and S¨umer, E. A computer vision system for chess game tracking. In *21st Computer Vision Winter Workshop, Rimske Toplice, Slovenia* (2016).

[3] Masouris, A., and van Gemert, J. End-to-end chess recognition. *arXiv preprint arXiv:2310.04086* (2023).

[4] Neufeld, J. E., and Hall, T. S. Probabilistic location of a populated chessboard using computer vision. In *2010 53rd IEEE International Midwest Symposium on Circuits and Systems* (2010), pp. 616–619.

[5] W¨olflein, G., and Arandjelovi´c, O. Determining chess game state from an image. *Journal of Imaging 7, 6* (June 2021), 94.