

# Final Report

**Chessboard State Prediction with Multitask Learning**

**Calum Wallbridge**

**Submitted in accordance with the requirements for the degree of  
Computer Science**

**2021/22**

**<Module code and name>**

The candidate confirms that the following have been submitted.

Items	Format	Recipient(s) and Date
Final Report	PDF file	Uploaded to Minerva (DD/MM/YY)
<Example> Scanned participant consent forms	PDF file / file archive	Uploaded to Minerva (DD/MM/YY)
<Example> Link to online code repository	URL	Sent to supervisor and assessor (DD/MM/YY)
<Example> User manuals	PDF file	Sent to client and supervisor (DD/MM/YY)

The candidate confirms that the work submitted is their own and the appropriate credit has been given where reference has been made to the work of others.

I understand that failure to attribute material which is obtained from another source may be considered as plagiarism.

(Signature of Student) \_\_\_\_\_

**Summary**

<Concise statement of the problem you intended to solve and main achievements (no more than one A4 page)>

Immediately explain what you've added compared to other implementations. beats the strongest openly available models.

### Acknowledgements

<The page should contain any acknowledgements to those who have assisted with your work. Where you have worked as part of a team, you should, where appropriate, reference to any contribution made by other to the project.>

Note that it is not acceptable to solicit assistance on ‘proof reading’ which is defined as the “the systematic checking and identification of errors in spelling, punctuation, grammar and sentence construction, formatting and layout in the test”; see

[https://www.leeds.ac.uk/secretariat/documents/proof\\_reading\\_policy.pdf](https://www.leeds.ac.uk/secretariat/documents/proof_reading_policy.pdf)

# Contents

<b>1</b>	<b>Introduction and Background Research</b>	<b>2</b>
1.1	Introduction . . . . .	2
1.2	Literature Review . . . . .	2
1.2.1	A <i>Brief</i> History of Computer Vision . . . . .	2
1.2.2	Computer Vision for Chess . . . . .	3
1.2.3	Prior Work From the Author . . . . .	5
<b>2</b>	<b>Methods</b>	<b>6</b>
2.1	Data Collection . . . . .	6
2.1.1	Sensors . . . . .	6
2.1.2	Auto-Labelling . . . . .	6
2.1.3	Dataset Versioning . . . . .	7
2.2	Model Training . . . . .	9
2.2.1	Experiment Tracking . . . . .	9
2.2.2	Board Segmentation . . . . .	9
2.2.3	Piece Recognition . . . . .	10
2.2.4	The Final Model . . . . .	14
2.3	Recording a Chess Game to PGN . . . . .	14
2.3.1	Leveraging a Chess Engine . . . . .	15
2.3.2	Motion Detection . . . . .	16
2.3.3	Towards Continual Learning . . . . .	16
<b>3</b>	<b>Results</b>	<b>17</b>
3.1	Piece Recognition . . . . .	17
3.1.1	Trials . . . . .	18
3.2	Recording PGN . . . . .	18
<b>4</b>	<b>Discussion</b>	<b>23</b>
4.1	Results . . . . .	23
4.2	Board Segmentation . . . . .	23
4.2.1	Piece Recognition . . . . .	23
4.3	Inference . . . . .	23
4.4	Dataset . . . . .	23

<b>CONTENTS</b>	<b>1</b>
4.5 Conclusion . . . . .	24
4.6 Ideas for future work . . . . .	24
<b>References</b>	<b>25</b>
<b>Appendices</b>	<b>26</b>
<b>A Self-appraisal</b>	<b>26</b>
A.1 Critical self-evaluation . . . . .	26
A.2 Personal reflection and lessons learned . . . . .	26
A.3 Legal, social, ethical and professional issues . . . . .	26
A.3.1 Legal issues . . . . .	26
A.3.2 Social issues . . . . .	26
A.3.3 Ethical issues . . . . .	27
A.3.4 Professional issues . . . . .	27
<b>B External Material</b>	<b>28</b>

# Chapter 1

## Introduction and Background Research

### 1.1 Introduction

Algorithms such as Deep Blue [4], AlphaZero [14] and more recently Player of Games[13] have enabled computers to out smart the smartest humans at the game of Chess. Unfortunately all these algorithms are bound to the digital world, rendered useless when competing against humans on a real board. This project aims to explore a major component of this: vision.

Unlike humans, the hard part of chess for computer is not planning which move to take next, but instead recognition, localisation and manipulation of objects in 3D space which currently all present much greater challenges. Perhaps the reason for the vision problem feeling so apparently effortless to humans is that over half of the human cortex is allocated to visual processing [15]. It is also described by Szeliski as an inverse problem and to attribute that as the reason for it's difficulty within computer science [16].

Consider the vision problem for chess to be two-fold: what is the current board state and where are all of the pieces? In particular this project will focus on the former, that is, to produce and present a solution for determining the state of a chess board from a video stream. A solution reliable enough to live up to the likes of AlphaZero in a robotic system, but also a solution that could be immediately useful in other applications such as realtime chess analysis from a real board.

There will be a focus on deep learning techniques, with consideration for best practice and the aim to share the tools to more easily manage and create new datasets in this area. Something called for by [8] as a serious challenge and priority for future research.

### 1.2 Literature Review

#### 1.2.1 A Brief History of Computer Vision

Computer vision is the study of making sense from visual data. Applications include character recognition for digitising documents, segmenting and classifying object instances for cancer screening, object tracking for sports analysis from video and countless more. To reach these goals, we name *features* to be useful pieces of information within an image (or any other visual data) and *feature detection* to be the class of algorithms that can extract features from images.

Some of the earliest work in computer vision started with Larry Robert in his 1965 paper [] describing a simple 2x2 convolutional kernel for edge detection which soon became the predecessor to the Sobel operator in 1968 [] and the still widely used Canny edge detector developed in 1986 []. The ability to find edges provided a helpful backbone for detection of higher order features such as lines and corners.

Another one of these tools that is still commonly used to help more sophisticated feature extractors is *thresholding* which refers to the process of partitioning images into two sub regions based on a pixel color/intensity threshold value. More complex partitioning schemes are often referred to as image segmentation techniques. Otsu's method is a famous example of an automatic global thresholding algorithm that finds a suitable threshold value for the entire image without supervision by minimising the resulting variance both between and within the two separated sub regions [1].

For higher order features more sophisticated methods are needed, a bucket term used to describe some of these methods is *template matching*, which can generally be considered as comparing an image to known feature templates to determine if that feature is present<sup>1</sup>. The Hough transform is the one of these higher level feature extractors commonly used to extract lines, circles and even arbitrary shapes, which heavily relies of edge detection [1].

Feature Descriptors (keypoint, interest detectors) ... Laplacian, Harris, SIFT, SURF, ORB, HOG

## Neural Networks

Enter... Neural Networks now do all of the above, better. Brief overview of some of the big ones then focus in on a few specific architectures for image classification used throughout the project

### 1.2.2 Computer Vision for Chess

Despite chess being a very narrow application of computer vision, the amount of research effort gone into the problem of determining board state is not insignificant. A variety of approaches have been tried and tested for which the following section will now summarise.

As in [8] the vision problem can be split into two further problems for analysis: board detection and piece recognition.

#### Board Detection

The problem of board detection is not specific to chess but also receives heavy research from other applications such as camera calibration [7]. The built in camera calibration functions in opencv [3] and matlab [12] are used in many previous works [11, 2] which provide a quick and precise solution for board detection, but becomes unusable when any artifacts like chess pieces are present on the board. This forced those authors to take the approach of an initial setup stage at inference, making the solution unfit to changes in board position during inference.

Due to a chessboard's simple features many early works of line and corner point detection can be applied. For example Hough transforms are used to detect the lines of a chessboard [10, 5]. Corner point detection methods such as the Harris and Stephens's [1] were also common among solutions [1], with some authors combining approaches with further processing such as canny edge detection [1] to yield more reliable results.

ChESS was another corner detection algorithm that out performed the Harris and Stephen's algorithm [1]. This was, perhaps interestingly, created for real-time measurement of lung

---

<sup>1</sup>Neural networks perform template matching on an input against their learned internal features

function in humans, further demonstrating the attention chess board detection has received due to its general applicability.

There are many other algorithms that require simplifications such as custom green and red chessboards [], multiple camera angles [], or even the requirement of user input for entering the corners of the chessboard [].

The most impressive work came out of Poznan University of Technology which proposes many interesting ideas that perform more reliably in a wider range of difficult situation such as pieces being present on the board []. They employ an iterative approach with each iteration containing 3 sub-tasks: line segment detection, lattice point search and chessboard position search. In each iteration of line detection a canny lines detector [] is used on many preprocessed variations of the input image to maximize the number of relevant line detections which are then merged using a linking function. The lattice point search starts with the intersection of all merged lines as input, converting these intersection points to a 21x21 pixel binary image of the surrounding area and runs them through a classifier to remove outliers. The addition of a neural network as a classifier greatly improves the generality of the proposed solution as it can be resistant to lattice points that are partially covered by a chess piece. The final sub-task then creates a heatmap over the original image representing the likelihood of a chessboard being present. Under the hood this is done by calculating a polyscore for the set of most likely quadrilaterals formed by the lines of the first stage. The polyscore is a function of the area of the quadrilateral and the lattice points contained within it. It is the quadrilateral that produces the highest polyscore that is used to segment the image for input to the next iteration until the quadrilateral points converge. The main disadvantage of this approach compared to others is that it can take up to 5 seconds to process one frame, for most use cases however this will be sufficient, unless realtime board training is required.

## Piece Recognition

Piece recognition has proved more difficult []. Most chess vision systems avoid classifying pieces by type (knight, king, ect.) all together []. These approaches typically get around this by requiring the board to start in a known position. From this known state the normal rules of chess can be used to infer what pieces are where after each move. Simpler methods require human input to prompt when a move has been taken [], more sophisticated attempt to do this move detection automatically.

These automatic move detection methods tend to all follow the same overarching processes of thresholding to detect a move having occurred, whether on color [] or even the edges detected within each square []. Most authors recognise the dependance this approach has on lighting variations, with Otsu thresholding [] sometimes being used to minimise the negative impact when lighting changed. While this improved results for what may be considered normal lighting conditions, they still suffered. They calculate reference colors for all 4 variations (white square, black square, white piece, black piece). All of these then only work in situation where a series of moves are to be recorded according to typical chess rules and not the chessboard state at any given moment.

There were a couple of methods that stood out from the rest each in their own way. One used

fourier descriptors to model piece types from a training set and the other modelled the pieces in Blender, a 3D modelling software, utilising template matching to determine piece type. The fourier method was very sensitive to change in camera angles, preferring a side view angle that unfortunately cause too many occlusions to be practical. The template matching approach took over 150 seconds on average to predict board state from one image which does not lend itself to interactive play in a robotic environment.

Go to standford dude and the heatmap guys as the best approach out there. They use SIFT and hard coded color algorithems. heatmap guys improved on this only by adding more restrictions by assuming the board much be valid and making statistical assumptions on what state is most likely. Oh and a HOG method. The one that said SIFT didn't work well because the lack of texture.

More recently another group of methods have surfaced using neural networks, specifically convolutional neural networks (CNNs) [1]. One of these used a pretrained Inception-ResNet-v2 model [2] and only had 6 classes, resorting to the more tradition approaches for color detection, in particular binary thresholding with added morphological transformations to reduce noise as seen in previous works [3]. Interestingly the six chosen classes were 'empty', 'pawn', 'knight', 'bishop', 'rook' and 'king\_or\_queen' as they claim kings and queens can be difficult even for human eyes to distinguish. Because of the choice of classes this method falls back to relying on a chess engine to determine piece type, which while usually correct for normal games of play makes the method unusable for games played with a variation on the normal rules of play. The other two methods used a simpler CNN structure similar to that of VGG [4] with 13 classes, one for every colored piece as well as the empty square.

### 1.2.3 Prior Work From the Author

Mention robotic arm and vision system for two counter board games as well as automatic differentiation library.

# Chapter 2

## Methods

To summarise, the overall approach for determining board state is to first segment each square of the board from the image, and use those subimages as input to a piece classifier from which the output can be concatenated into a view of the whole board. Splitting the input to a classifier up by square is a unanimous decision taken by previous authors discussed in Literature Review as it reduces the problem to a simple classification problem of only piece type and can take advantage of the thoroughly researched area of chessboard detection.

### 2.1 Data Collection

At the heart of any machine learning project is the data. It is as important, often more important, than the code and presents many interesting challenges. **Why is this?\*\*** Discussed in the following sections are some of the challenges and decisions that were considered.

#### 2.1.1 Sensors

The eminent challenge is acquiring data in the first place and is highly context dependant, For vision there are a range of sensors we can use to gather data from the real world. Sensor choice is an important choice for any robotics application as there are important tradeoffs, as with any engineering challenge, which must be considered.

**Outline some of the tradeoffs between spatial sensors\*\***

One important distinction to make is the difference between training and inference. Requirements at the time of training my differ significantly to the requirements at inference. Processing power, energy supply and realtime operation are some of the constraints that will have to be met when considering different sensors.

Talk about single camera.

The sensor used throughout this project is the RealSense SR305 which is a RGB-D camera using structured and coded light to determine depth, it functions best indoors or in a controlled lighting situation. For the reasons outlined above the RGB camera stream is mainly relied upon but there will be some discussion and comparision of piece detection with the depth sensor.

Talk about the generic camera

#### 2.1.2 Auto-Labelling

Talk about using a simulator.

A closely related challenge of acquiring the data is that of labelling it too. During Literature Review multiple past authors have stated the availability of datasets for chess piece recognition

is sparse [] with some emphasizing dataset collection took the large majority of their time []. It is also widely known that neural networks scale with the number of examples [], which will be explicitly explored for Chess Vision in Results. This however poses the question: how do we get access to a lot of labelled data for chess?

Unlike techniques in [] \*\*Some examples of other auto labelling techniques\*\* the approach taken here was to maximize speed of collection and flexibility.

Portable Game Notation (PGN) is a common format for recording chess games as a series of moves and are widely available online. All the PGN files used in the project were used from [] which has over one million games. Utilising this data not only has the benefit of an abundance of chess games but also that the recorded games are real and contain positions more likely to appear in game play.

A program is developed for recording these games with the generic camera interface as previously described. The program takes screenshots upon user input (with the [Enter] shortcut) displaying the move number and image as a result for visual feedback before saving to disk. These games can then be automatically labelled using the matching PGN file.

After the development of this pipeline it was possible to collect over 2,500 *unique* labelled images in under two hours.

### 2.1.3 Dataset Versioning

With all this data the next challenge becomes self evident. It is concerned with the question: How do we manage all of this? As experiments are carried out and iterations on the dataset are performed there will be many changes and variations of the data that are used to train models. This is a challenge with data for many reasons, the first being reproducibility. Say you train a model that performs really well, then you make some changes in say the balancing of the dataset or even the size and you train again. If you go through a series of changes like this and discover that it performed a lot better the way you had it previously then you might want to roll back.

In typical software engineering a version control system like git [] would be used. While this is what was used to version the codebase of this project, there are different challenges with data, namely it's storage. Code usually takes up megabytes of storage, the Linux Kernel source code for example contains 27.8 million lines and is only around 1GB in size []. The data used throughout this project came to >20GB.

There are many proposed solutions for this problem, from managed service like Neptune [] and wandb [], to self hosted opensource options. In fact Git has its own solution called Git Large File Storage (Git LFS) [] which uses the exact same methodology as with code except it will store large files in an external remote and only reference those locations in your git repo using a hash of your data's content. This means if you change any of your data a new copy of that data will be stored and it's reference updated. And so to maintain git's methodology of versioning any change of your data will mean a new copy will be stored, rapidly increasing storage costs with changes. Neptune instead is a more holistic machine learning platform that provides a whole bunch of features other than just dataset versioning. As with other managed

services these managed solutions are prone to lock-in and require you modifying your code with a bunch of API calls to get them running. In the end it was found that a custom solution utilising some cloud based object storage proved most useful, with the least amount of effort and cost. Perhaps the wide variety of solutions all with different approaches is evidence of this not yet being a solved problem. To demonstrate the cost difference, Git LFS per GB cost is \$0.1 a month which is 4 times more expensive than amazons most expensive rate of \$0.023 per GB. With optimisations you can get the amazon S3 bucket price down to \$0.0125 making Git LFS 8 times more expensive. These numbers may look small but they add up with time and scale, especially when versioning many changes where copies and diffs need to also be stored.

The solution used in this projects centers around 3 elements: The *Game*, a *Labeller* and the *Storage* facade.

The storage facade's purpose is to abstract file storage so that it could work with any backend and provide a simple to use interface for fetching files using a file system.

`file = open(Storage("img.png"))` for example will give you the file descriptor for img.png cached from the filesystem if it exists, pulling it from an external store if not. This was very useful for training across different VMs in the cloud and could work with any storage implementation. The only backend implemented was Amazon S3 bucket store, and used less than 50 lines of code.

The Game class ?? is how chess data specifically is managed. Each instance of Game represents one 'unit' of data where each unit is a recorded chess game. You can record a game by using the recorder application described in the autolabelling section and is simply a sequence of images with an associating PGN description. *LabelOptions* can be set on a game which will be used by a Labeller to describe exactly how the autolabeller should function. For example, `Game("Kasparov", LabelOptions(margin=50))` will create a game that should be labelled with a margin of 50px surrounding each square. As will be discussed later it is this game object that will version the data used to train a model while dramatically decreasing the storage cost. This is opposed to storing a new entire labelled dataset everytime a slight change is made, i.e chaning the margin of each square.

As different model architectures were explored, different labelled data was need entirely. For example, one model may only be predicting whether a piece is a black or white pieces whereas another may be prediting the type of piece. To cater for this variance while keeping the rest of the training pipeline unchanged the Labeller abstraction is used. Each labeller has a function that can take in a Game and output a series of X, Y input, target pairs. These can then be saved to disk for the *ChessFolder* pytorch dataset to handle. This could even be extended to more complicated labels like bounding boxes or multiple input images.

One important note in implementation is the use of python generators which dramatically reduce memory usage and speed up the auto-labelling pipeline by over 10 fold in places.

## 2.2 Model Training

### 2.2.1 Experiment Tracking

Throughout this project, over 2000 models were trained, each slightly different, whether in architecture, hyperparameters or the data used. This makes the processes which is commonly referred to as Experiment Tracking a very important one. It would be good to visualise all these experiments and be able to roll back to old models. As with dataset versioning there are a plethora of new managed services that promise they can do this for you. The approach opted for during this project utilised an open-source solution called Guild [1]. The main reason for this choice was because of its unopinionated nature, requiring zero code changes. GuildAI leverages the filesystem to save *runs* including the environment configuration at the time of training. A run is made from the command line (i.e.

`guild run train --remote ec2 EPOCHS=10 LR=0.001`) and represents one iteration of a model and can even be executed remotely via ssh to give easy access to GPU machines. It can also perform hyperparameter optimization using built-in techniques like grid search, random search, gradient boosted regression trees or even your own custom optimizer. To enable the zero code change promise, guild fetches hyperparameters from a variety of places including globals set within your training script, configuration files and commandline arguments. This method was greatly preferred to using a paid managed service or some other solution which requires a database to setup. The previously mentioned dataset versioning method also greatly benefits from guild's use of the filesystem as within each run we can also store the Games and Labeler that was used to train the model. TensorBoard integration also played a huge role as it provided an easy to use, visual place to interpret and compare runs. See Figure 2.1 for examples.

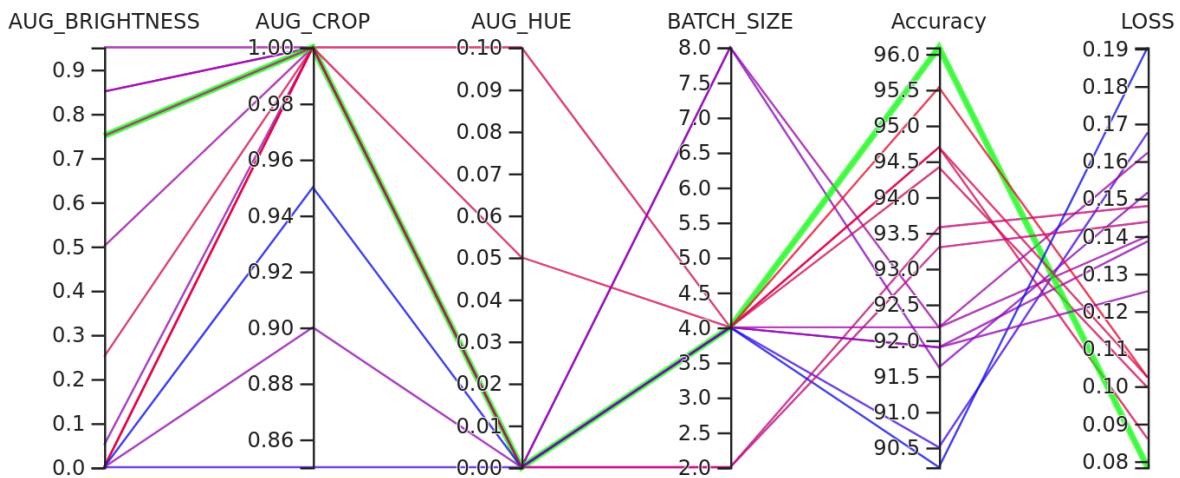


Figure 2.1: A parallel coordinate view of runs within TensorBoard

### 2.2.2 Board Segmentation

Aruco markers are chosen for board corner point detection as a very simple method that works reliably even with a board full of pieces. With the corner points of the board the perspective transformation is calculated using Gaussian elimination [?], as demonstrated in Figure 2.2.

Although Aruco markers require customizing the environment they are very fast at inference and allow the focus to remain on piece recognition, from which Literature Review showed is less studied and reliable. In a more holistic solution the iterative heatmap approach proposed in [6] would be recommended.

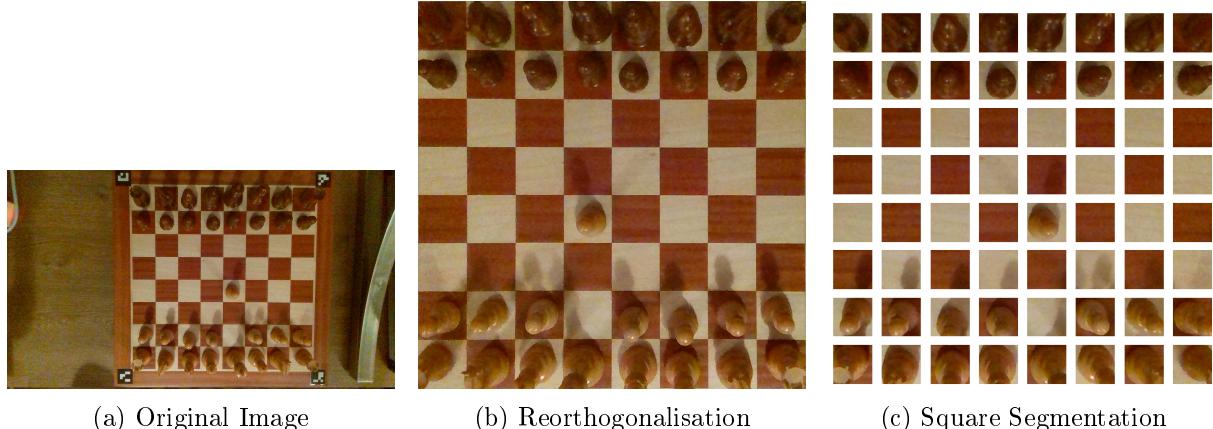


Figure 2.2: Square segmentation process. The 64 images in (c), after augmentation, is what ultimately gets sent to the model for classification

### 2.2.3 Piece Recognition

Now that the board has been segmented including all of it's squares, it is time for the fun stuff. That is to determine what piece, if any, occupy each square. To start, a good baseline is found. A good baseline is a simple model to understand and easy to get decent results with. For classification, the pathological baseline would be a uniformly random model, which could easily be extended to use a categorical distribution. The probability mass function for a categorical distribution of  $k$  categories numerically labelled  $0, \dots, k$  is

$$f(x \mid \mathbf{p}) = \prod_{i=0}^k p_i^{[x=i]}$$

where  $\mathbf{p} = (p_0, \dots, p_k)$  and  $p_i$  represents the probability of an image of category  $k$  being sampled from the training set. Since  $\sum_k p_k = 1$  you can generate the probabilities by normalising the count of each category in the training set. Algorithmically sampling from the distribution can be done using inversion sampling which requires calculating the cumulative distribution function.

Starting with a known baseline is common practice in exploratory machine learning [] so that transitions are always made from a known and working state. It becomes very easy to see if an experiment is not working by comparing it to your baseline.

Keeping to this strategy, the Multilayer Perception (MLP) or fully connected network [] was the first neural network to be explored. We explore nerual network approaches as they have proved most effective for image classification [], they are also relatively simple to deal with as no hand-crafted feature extraction is needed and end-to-end solutions are much more viable. By starting with the MLP, all complexity from the network is stripped away so the more

extraneous elements such as the training loop and evaluation metrics can be built and tested.

Evaluation metrics were very important in developing as loss is not that human interpretable. Firstly accuracy, which is much more interpretable, is calculated to be the ratio of correct classifications to total number of sample images. This could then be extended to top-n accuracy which takes correct classifications (true positives) to be the count of classifications where the actual class was within the top  $n$  confidence scores of the models output. This is useful for problems where there are a large number of classes or when at inference there is the intent to do a search through the output probability distribution which is more relevant to this project as there is at most only 13 classes and a search through the probability space at the inference stage is entirely possible.

As the number of evaluation metrics grew it became prudent to encapsulate them into an *Interpreter* class so as to not clog up our training loop which itself is within a *Trainer* class to encapsulate its implementation away from the training script which will see frequent changes during experimentation. Some of the first additional methods added to the Interpreter is `plot_top_losses` and `plot_confusion_matrix` to make it significantly easier to find which classes the model was confusing and which specific images it found difficult to classify giving hugely helpful insight into the dataset itself and even finding bugs in the autolabelling pipeline.

Before training, a fixed random seed was important so that it was possible to tell whether or not any particular change was positive since all models started from the same initialisation. The initialisation of the final layer was hard coded to be an equal distribution, sensible initialisations of final layer can have a noticeable effect on training convergence [] .

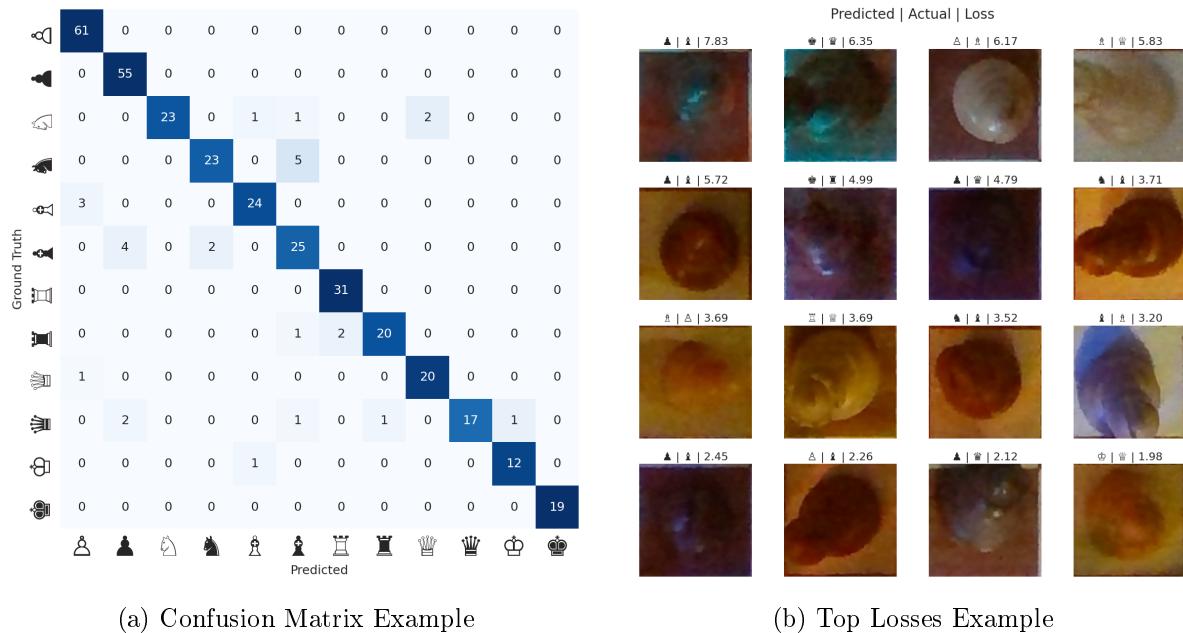


Figure 2.3: Plots that are saved for each run and can be viewed in TensorBoard

Once the full training and evaluation structure is functional, new features can be incrementally added and architectures explored. - human labelled dataset?

Another strategy was to purposely overfit models during training. It was easy to tell when a model was overfitting by splitting the data into training and validation sets [] and viewing the

relative losses in TensorBoard as can be seen in ?? . The reason for doing this was that if the model was able to overfit (i.e. achieve  $\sim 100\%$  accuracy on the training set) then it was able and big enough to learn distinguishing features or as is sometimes put: the model has sufficient entropic capacity. It is then a lot easier to make it generalise well, for example by reducing the number of parameters, than to go from an underfitting model to a generalised one as you don't necessary know where the problem lies for not fitting to the data.

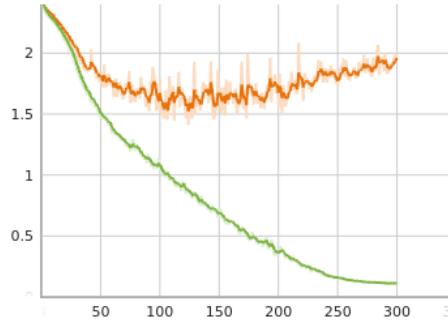


Figure 2.4: Demonstration of overfitting. Loss against epoch. The validation loss in orange is seen diverging above and the training loss in green is seen tending to zero below<sup>2</sup>. A model should be saved at roughly epoch 125.

There are many ways to fight overfitting as typically referred to as regularization techniques, perhaps the easiest of which is early stopping []. The way this was implemented in our Trainer was to save the model after the first epoch and consequently after every other epoch for which the validation loss was less than the previously saved model. It is possible to then stop training entirely if after a set number of iterations no improvement is seen in the validation loss.

As mentioned in Literature Review convolution operations, and in particular differentiable convolutional operations have had monumental impact on the field of computer vision and so this was the next experiment. Quite quickly, especially with a limited dataset, overfitting became a major problem to overcome and so many experiments were positioned to solve this problem. Below are an overview of some of those experiments including some final optimisations to squeeze as much performance out of our model as possible.

## Architecture

ConvNext <https://arxiv.org/pdf/2201.03545.pdf>

Due to the recent surge of deep learning research for computer vision there are a huge range of openly available models. During this project ResNet, ConvNext, and ViTs were experimented with. How much detail should I go into here?

ResNet

## Loss Function

<https://towardsdatascience.com/choosing-and-customizing-loss-functions-for-image-processing-a0e4bf665b0a> <https://arxiv.org/pdf/2009.13935.pdf> proposes SML may be better than cross-entropy

## Optimizer

[put into results] AdamW was found produce more attractive results from comparing SGD with/without momentum and Adam and AdamW [] Learning Rate, Weight Regularization. Learning rate scheduler, leave this till the end.

## Batch Normalization

<https://arxiv.org/abs/1502.03167> and renorm for small batch sizes

<https://arxiv.org/abs/1702.03275>

<https://arxiv.org/abs/1207.0580> shouldn't be used after convolution layers [] there has been some work <https://arxiv.org/abs/1904.03392>

Pooling and DropPaths

## Transfer Learning

There appears to be a trend occurring in the deep learning space. Some organisation spends millions training an impossibly large nerual network and others more and more are using these models, often fine-tuning for their own use cases. [] uses these large models as fixed feature extractors.

This appraoch makes sense as it is impractical to retrain huge nerual networks that take weeks, millions of dollars and wasteful amounts of energy to train [].

In the case of CNNs we can see the features that kernals in the early layers learn [] are often very simple shapes and will be common for all computer vision tasks. This will explored this further in the results section as we visualise kernals from both random initalised models and pretrained models.

## Ensemble and Multitask Learning

Another regularization technique is known as ensemble learning and is common practice in deep learning to reduce variance inherent to stochastic training with random initalisation [9]. The technique involves training several models (with different random seeds or configurations) and combining their outputs. Using an ensemble almost always beats training a single model as two models are unlikely to make exactly the same errors on an unseen dataset. This projects extends on this idea to also combine models with different class predictions.

In the two previous deep learning solutions for chess piece recognition, neither researched into the effect of using differnt class groups, despite chosing differently. In this work, 5 sets of classes were considered as shown in Figure 2.5. Models trained on each set of classes were compared with accross each metric and some were selected to be used together. For example, combining the predictions of 3 model separately trained with the 'piece', 'color' and 'occupied' labellers would yeild the same class of predictions as in 'all'.

Labeller	Count	Classes
all	13	Empty, White Pawn, White Rook, ..., Black Queen, Black King
piece	12	White Pawn, White Rook, ..., Black Queen, Black King
occupied	2	Empty, Occupied
color	2	White, Black
type	6	Pawn, Rook, Knight, Bishop, Queen, King
type+	7	Empty, Pawn, Rook, Knight, Bishop, Queen, King

Figure 2.5: Class Sets

Every model you add to an ensemble has a large impact on the computational cost at inference (as well as training) which must be considered. A separate approach called multitask learning which separates heads for different classes in a network might be worth considering instead.

Part of the benefit of separating models for different classes (which can be referred to as tasks in this context) is that they each get their own parameters they can optimise. By using 'all', each task has to fight for optimisation. Some tasks may train well together and benefit from sharing parameters, some may not. Multitask learning is the study of finding this balance of which parameters should be shared among each task.

It is also worth mentioning loss functions here as they will need some attention if multitask learning is to be used. Typically a separate loss function will be used for each task and then these are simply summed together before performing backpropagation. One useful technique is to be able to weight these losses towards the more difficult tasks. This weight vector then becomes another hyperparameter to optimise. It's good to know that this technique can be used with many individual loss functions across specific classes too.

## Augmentation

Data augmentation is a strategy every machine learning practitioner has in their arsenal. It addresses the data problem, allows us to truly leverage the data we have and generalise our models further. The MNIST [] dataset itself was created using data augmentation. In the case of chess piece classification, the correct augmentations or transformations should be chosen. Go onto list the transformations used and why.

### 2.2.4 The Final Model

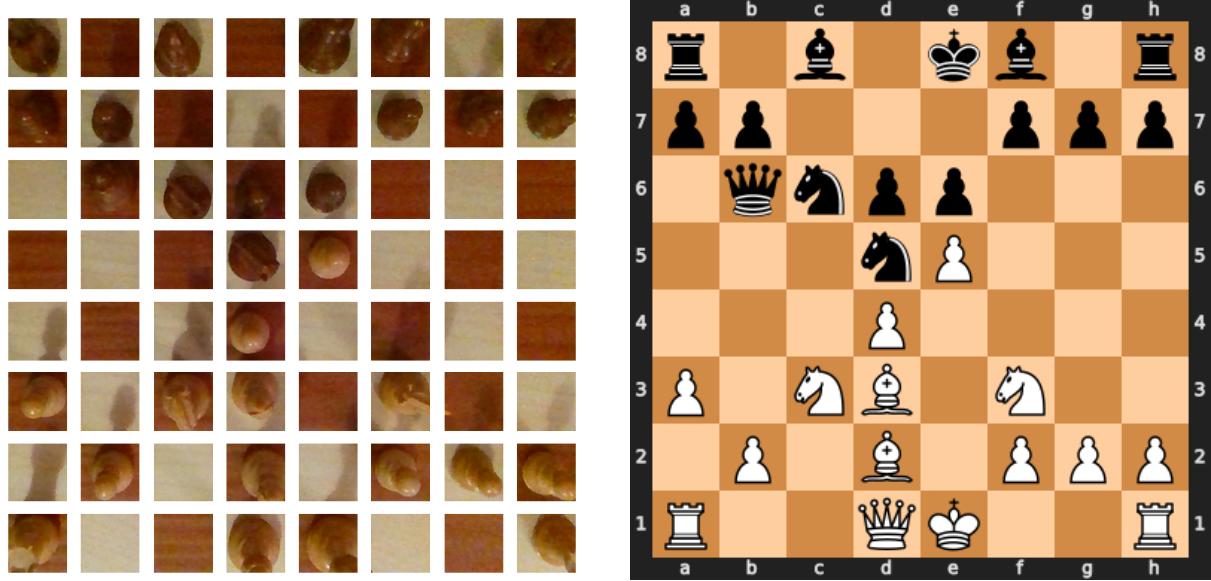
Augmentation as well? Talk about grayscale (a lot of previous works used gray scale do performant CNNs aid from gray or hinder?)

## 2.3 Recording a Chess Game to PGN

The goal of this section is to discuss methods for using the proposed piece classifier to record an entire game of chess, played on a real board, to a PGN formatted file. This application is later referred to as the *Inference Application*.

The general approach is to generate a board state at each fetched frame. That is to segment the board and each of its squares, send each square through a forward pass of the piece

classifier described above and finally collate the predictions together into a board state. A board state is a generic term that in actuality could be many things, but in this case it is sufficient enough to think of it as Forsyth-Edwards Notation (FEN). This board state can then be compared with the previous board state and if any difference is detected then that move is added as a child node. This tree structure gives flexibility for many variations of the same game to be recorded to later be parsed and encoded to PGN.



(a) Frame after Segmentation

(b) Prediction Vision State

Figure 2.6: Board Squares Segmentation Process

### 2.3.1 Leveraging a Chess Engine

As leveraged by others before [1] a chess engine and other statistical methods [2] can be used to increase the reliability of board state prediction when the normal rules of chess can be assumed to be abided by. The proposed piece classifier in 2.2.4 makes no such assumptions, which lends it self to a wider array of circumstances even when the normal rules of chess are not being used. When a game is to be recorded to a PGN file however there are some rules you must follow so that a game can be encoded. Typically chess engines can handle a few variants and so it is assumed to be safe to incorporate one here without loss of too much generality. And so a chess engine is used, not to help the classifier's predictions, but to automatically determine when a move has been made. Specifically leveraging a chess engine means that we assume (or set) the starting state of the board and only allow legal moves within a chosen chess variant to be accepted, much like a human would when playing against a competitor.

To do this, two concepts are introduced: *VisionState* and *BoardState*. With some simplification, the VisionState is a lightweight representation of board as the model sees it in the last fetched frame. The BoardState starts at the known starting state and is only updated if the VisionState at any time step represents a legal move from the current BoardState. When the BoardState reaches a terminal state, or otherwise receives a cancel signal in cases of a draw or resignation, the game tree is parsed and the PGN saved to disk.

Since our classifier is almost guaranteed to never achieve 100% accuracy on new unseen data, especially as more chess sets of different shapes and sizes are used, it is likely to have varying uncertainty. Often, the more uncertain predictions of a particular square may result in flickering, that is for example it deduces the piece to be white in one frame yet black in the next. This occurs when the resulting probability distribution for a given input looks something like  $p = [0.38, 0.39, 0.09, 0.04, 0.1]$  as it is not unlikely with the next frame we see it change to  $p = [0.39, 0.38, 0.09, 0.04, 0.1]$  resulting in a different classification. To cater for this uncertainty like this, the VisionState also has a concept of memory with length  $N$ , where memory is in an average of states over the last  $N$  frames.

### 2.3.2 Motion Detection

One factor that was found to still sometimes break this system was motion. Moving pieces across squares and hands flailing over the board confused the model. Even with the separation of the BoardState and memory to remove anomalous predictions - especially when motion persisted for longer periods.

In a lot of these situations the actual board state is undefined. That is because a piece has been lifted and so must now be moved but not yet let go and so its definition may be undetermined. Because of this it is not unreasonable to halt inference all together. User input to indicate when a move has been completed is a common strategy [], but goes against the purpose of building such a system all together - autonomy. Instead a motion detector is employed. Even the naive motion detector of using a threshold over the absolute difference between each consecutive frame was found to be sufficient for removing these disturbances. Some other methods such as SIFT and SURF were also explored but found to be more computationally expensive than necessary.

### 2.3.3 Towards Continual Learning

Another feature added to the Inference Application was the ability to save snapshots. During inference it was not uncommon to see the vision state make mistakes, and while the above features usually enabled us to generate an accurate PGN file, this is still useful data. By giving the user the ability to save snapshots when they see an error it is now possible to feed this data back into fine-tuning the model. This is the process of continual learning and is an active area of research [].

# Chapter 3

## Results

This chapter will cover the results of the best and final model that was trained. The first section, Piece Recognition, measures the performance of piece recognition against existing solutions. The second section, Recording PGN, tests the model within the inference application to provide a better view of its utility.

### 3.1 Piece Recognition

The best openly available solutions for chess piece recognition on real boards were found to be from NAME [] and NAME []. It was hypothesized that through each iteration and optimization of hyper-parameters the presented solution was being overfit to the images in the evaluation set despite never being directly trained on it giving the proposed solution an unfair advantage when tested on the original board. And so, in order to make a fairer comparison, another board entirely was also chosen for these final tests. Including a dataset from an unseen board should allow for more meaningful conclusions to be drawn regarding the generality of each model.

Three tests, with the results shown in Figure 3.1 were then performed. The same training and evaluation splits were used for each model and the metrics were averaged over 3 runs. The first test used data from the original board only, the second with data from only the unseen board and the third test included data from both boards.

[] Uses a pretrained VGG16, freezing all layers of the convolutional feature extractor, adding a head of 3 fully connected layers with a sum of 2.5M trainable parameters. The final layer outputs a softmax distribution over all 13 class ('all' from Figure 2.5). After getting familiar with the code base and running some different experiments, a few minor improvements were spotted that could be made without changing the architecture, such as implementing early stopping. In order to ensure a fair representation of the authors work those changes were implemented and the best results were used in Figure 3.1.

[] Trains a Support Vector Machine classifier for each piece type including empty ('type+' from Figure 2.5) using the HOG features extracted from each training image. For each input image the classifier with the highest confidence determines the piece type, to determine piece color they use reference colors. This involved manually calculating reference colors for each piece and deciding on a suitable threshold value which was done separately for each board.

Unfortunately, this method made it difficult to generalize when using multiple boards and so the decision was made to exclude it from the final test.

Weighted accuracy is the typical accuracy calculation:

$$Accuracy = \frac{\sum_{i=1}^k tp_i}{\sum_{i=1}^k tp_i + fp_i} \quad (3.1)$$

Method	Original Chessboard		Unseen Chessboard		Both Chessboards	
	Accuracy	Balanced	Accuracy	Balanced	Accuracy	Balanced
proposed	<b>0.97</b>	<b>0.95</b>	<b>0.96</b>	<b>0.96</b>	todo	todo
	0.87	0.71	0.86	0.75	todo	todo
	0.77	0.57	0.74	0.48	-	-

Figure 3.1: **Evaluation Accuracies** for models trained on both single chessboards and then a dataset containing images from both boards

$$\text{Balanced Accuracy} = \frac{\sum_{i=1}^k \frac{tp_i + tn_i}{tp_i + fn_i + fp_i + tn_i}}{k} \quad (3.2)$$

### 3.1.1 Trials

For a clearer comparison in the holistic task of full chessboard state recognition a similar methodology to that of Ding, Czyzqskil et al. [8, 6] is used. It involves collecting a small yet diverse benchmark of chessboard images and evaluating the full chessboard prediction from each model. This is in contrast to using the dataset used above which consists only of individual pieces and their corresponding labels.

Figure 3.3 demonstrates the extent of the challenge by visualizing the actual predictions from each model along with the input and ground truth. The results show... as in Figure 3.2 the overall performance.

Interestingly || method, while underperforming in accuracy on the evaluation tests, has a much higher recall of empty classifications compared to || who utilizes a CNN as in the proposed solution which beats even ||. It is proposed that the reason for this is because of the multitasking element, which is what was also found in the experimentation to find The Final Model which is not what was expected. What was expected is that sharing weights between the 'type' classifier and 'occupied' classifier would aid in identifying empty squares since the features that make up an empty square are just the lack of features that make any of the classes in the 'type' classifier. It turns out, however, that the empty classifier trains more effectively when it has its own weights it can adjust as the proposed model shares very few weights between the 'occupied' classifier and every other class and the empty classifier in || is completely independent as opposed to in || where all the features are shared.

Method	Mean Inference Time	Total Misclassified Squares	Cross Entropy Loss
Proposed	1.7s ± 61.4ms	-	3.01
	<b>1.5s ± 58.7ms</b>	-	3.21
	15s ± 417ms	-	1.82

Figure 3.2: Games played with the inference application

## 3.2 Recording PGN

In addition to evalutaing the raw model's performance, five games of varying length are played from beginning to end with the inference application. The games were played as normal with

the only exception being that moves were only made after the inference application detected the previous move. These delays, as they'll be referred to in Figure 3.4, were measured with a timer and recorded. If the model failed to correctly generate the correct PGN at the end of the game, or otherwise crashed, it was marked as a failure.

The one failed game (Carlsen/Nakamura) was due to the inference application registering a rook movement before the player lifted their hand from the rook who then decided to move it to a different square.

The fastest move to register took 0.27s and the slowest took 76s which appeared to be when there was lots of specular light reflecting from the pieces causing the model to confuse colors.

While the delays on the tail end of the distribution are certainly caused by uncertainty in the model, the cause of the median delay is more difficult to pinpoint and so a profiler was used quite a lot throughout the project and is shown here in Figure 3.6. A reasonable chunk of time is consumed by the display function which provides the user feedback as can be seen in the provided demo. As this is not vital to the function of the application it would be possible to place this logic within another thread with a lower priority. This should completely remove ( $\tilde{20}\%$  of the total delay) if the machine has a free CPU to schedule the thread to. Due to using a GPU and inference the entire time spent forward propagating the images through the model took less than 5% of the whole inference loop, the most glaring time during inference seems to be spent in the `.item()` method of pytorch which is moving data from the GPU to the CPU.

Input	Ground Truth	Proposed		
			Mistakes: 0	11
			Mistakes: 0	12
			Mistakes: 0	12
			Mistakes: 0	12
			Mistakes: 0	12
			Mistakes: 0	12

Figure 3.3: A sample of six segmented chessboard images and their corresponding state prediction from each model

Game	Total Moves	PGN 100% Correct	Delay for Move to Register (s)		
			Median	Mean	Std Dev
Carlsen/Wesley	57	✓	1.25	3.01	5.26
Carlsen/Rapport	36	✓	1.13	3.21	4.97
Carlsen/Nakamura	71	-	0.95	1.82	2.09
Carlsen/Aranian	55	✓	1.24	4.20	11.31
Fool's Mate	4	✓	0.95	1.19	0.76

Figure 3.4: Games played with the inference application

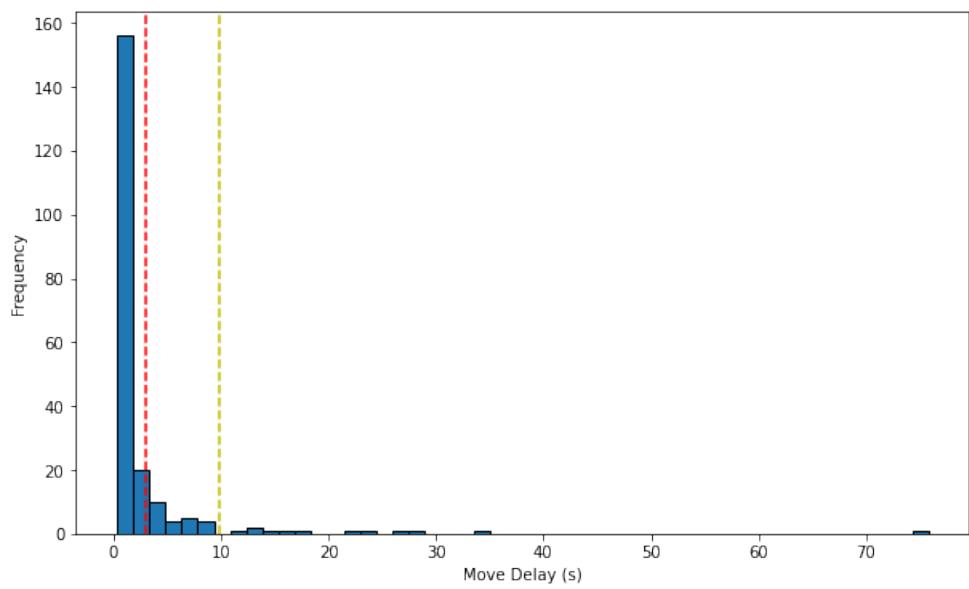


Figure 3.5: Histogram of move delay across all games in Figure 3.4 with arithmetic mean in red and sample standard deviation in yellow

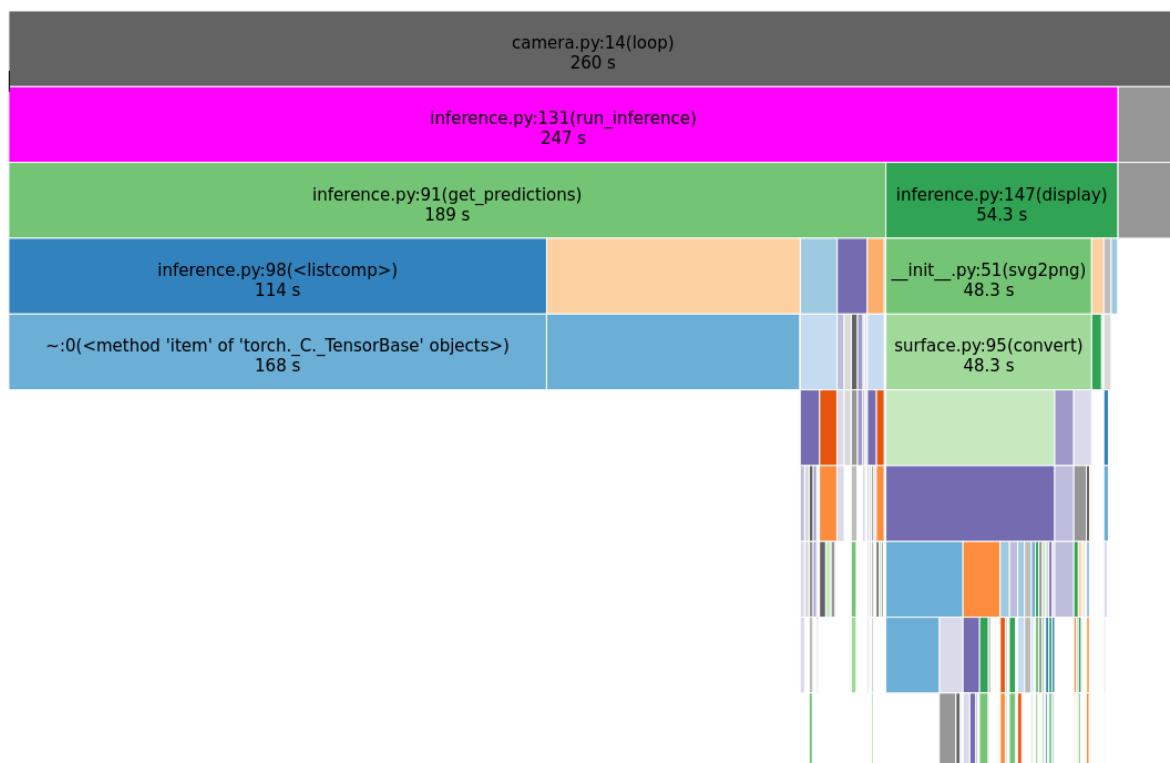


Figure 3.6: cProfile Visualisation of recording a 4 minute game at inference with moves made in quick succession

# Chapter 4

## Discussion

### 4.1 Results

### 4.2 Board Segmentation

The decision to use aruco markers, while impractical for release, aided development by being an accurate localizer. To correctly classify pieces, the segmentation of board squares was critical as a small error of even 5mm could throw the model off as due to added perspective shift a square could now contain two pieces. Compared to the iterative heatmap method [1] aruco marker are not the best, but in a controlled environment this did not become a problem and more importantly they worked even with pieces on the board which enabled recalibration during inference and while collecting data. This is unlike a lot of other proposed solutions [2]. Secondly they're fast [3]. The iterative heatmap method can take around 5 seconds to segment the board which while not terrible can add up when relabelling lots of data for many model experiments. As this project focused on piece recognition and the development of an inference system, aruco markers proved to be a sensible choice. However, they come with deal breaking consequences if this method is to be used in production with many boards. It's just too impractical to print and fix markers every time a new board is to be used with the system and again goes directly against the main purpose of this project: autonomy.

#### 4.2.1 Piece Recognition

The main component that appears to have contributed to the success of the proposed solution is the use of multitask learning combined with the state-of-the-art CNN architecture, ConvNeXt. Without exploring which classes benefitted from sharing parameters and the time put into hyperparameter tuning with the help of careful experiment tracking the results were a lot worse. Another important difference was the use of transfer learning, without parameter freezing like [4]. This contributed to an accuracy increase of ~10%.

### 4.3 Inference

Surprised at how well it worked.

### 4.4 Dataset

I found not only managing datasets but choosing and building them to be a big yet fun task. No doubt by the legends who've built tools for me to use, my experience has been that most of

the work is in the data. Creative ways to collect it, label it and fit it to pre-existing network architectures.

It was not until I recorded and labeled lots of data (1000s images) was it until I suddenly started to see new patterns I had not seen before. Small things like dark shadows, motion blur, pieces half across squares, poor resolution or exposed frames.

More control over the sensor hardware sensors may have been useful. Changes in the way we generate image (as in the autolabeller) or adding augmentation such as cropping and blurring helped which these changes that I would not have seen in my metrics before hand.

## 4.5 Conclusion

## 4.6 Ideas for future work

Firstly it would be nice to explore these methods with more extreme camera angles. This would probably include extending the labeller too add more margin in one direction to account for the perspective. As my model performs better than [] who found decent performance with other angles, mine should do pretty well, if not better.

If assumptions about the environment are to be kept minimal then localising pieces in 3 dimensional space should be a requirement for robotic manipulation to be possible.

Perhaps worth taking a step back and reconsidering the unanimously made decision by which chessboard state recognition is done. Instead of splitting the board up into squares and using simple image classification it would be more useful for robotic systems to have a 3D representation of the space. It is the author's belief that "3d reconstruction" would kill two birds with one stone and the direction that future research should focus.

Explain. Similar work in other areas.

# References

- [1] D. Ballard. Generalizing the hough transform to detect arbitrary shapes. *Pattern Recognition*, 13(2):111–122, 1981.
- [2] d. bowers. Robot arm, chess computer vision, Jul 2014.
- [3] G. Bradski. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, 2000.
- [4] M. Campbell, A. Hoane, and F. hsiung Hsu. Deep blue. *Artificial Intelligence*, 134(1):57–83, 2002.
- [5] W. L. K. Chua Huiyan, Le Vinh. Real-time tracking of board configuration and chess moves. National University of Singapore CS4243, 2007.
- [6] M. A. Czyzewski. An extremely efficient chess-board detection for non-trivial photos. *CoRR*, abs/1708.03898, 2017.
- [7] A. de la Escalera and J. Armingol. Automatic chessboard detection for intrinsic and extrinsic camera parameter calibration. *Sensors (Basel, Switzerland)*, 10:2027–44, 03 2010.
- [8] J. Ding. Chessvision : Chess board and piece recognition. 2016.
- [9] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016.  
<http://www.deeplearningbook.org>.
- [10] J. Hack and P. Ramakrishnan. Cvchess: Computer vision chess analytics. Stanford CS231A, 2014.
- [11] C. Koray and E. Sümer. A computer vision system for chess game tracking. 2016.
- [12] MATLAB. *version 7.10.0 (R2010a)*. The MathWorks Inc., Natick, Massachusetts, 2010.
- [13] M. Schmid, M. Moravcik, N. Burch, R. Kadlec, J. Davidson, K. Waugh, N. Bard, F. Timbers, M. Lanctot, Z. Holland, et al. Player of games. *arXiv preprint arXiv:2112.03178*, 2021.
- [14] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, et al. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, 2018.
- [15] R. Snowden, R. Snowden, P. Thompson, and T. Troscianko. *Basic Vision: An Introduction to Visual Perception*. OUP Oxford, 2012.
- [16] R. Szeliski. Computer vision algorithms and applications, 2011.

# Appendix A

## Self-appraisal

<This appendix should contain everything covered by the 'self-appraisal' criterion in the mark scheme. Although there is no length limit for this section, 2—4 pages will normally be sufficient. The format of this section is not prescribed, but you may like to organise your discussion into the following sections and subsections.>

### A.1 Critical self-evaluation

Firstly found it to be a really fun project and enjoyed writing all the components that make a machine learning solution. Seeing the end result was even better as I'm genuinely surprised at how well it worked, beating all the open source solutions I could find out there. Learnt a lot about statistics which I never covered at A-Level.

### A.2 Personal reflection and lessons learned

Surprised at how effective transfer learning is and the significance of it's place in the future. The importance of experiment tracking for research.

### A.3 Legal, social, ethical and professional issues

#### A.3.1 Legal issues

Legal issues could only really enter my project through the data collection and usage aspects. By using my own data of games I played with no personal information I was able to completely avoid any legal issues. As I am opening my data and solution to be opensource as I found that I really appreciated when others did so legal issues may also enter such as licensing. Any datasets collected with my tool will be the complete responsibility of the person who captured that information and would need to be included in the license.

#### A.3.2 Social issues

The one social issue I can really see with this project is if it is used for cheating. I also believe it would be trivial to train my solution on digital games such as those on chess.com as the data is abundant and would probably work very well. This could lead to some bad actors using it to run analysis while playing the game and so would ruin the game for those who want to play against real humans.

### A.3.3 Ethical issues

Personally, there weren't any serious ethical issues except those that I have already mentioned. The real ethics come into when the tools and techniques I have used throughout this project are used for other purposes. Deep learning is improving at a rapid rate and because of it's nature it's applicability knows really no bounds. I believe these techniques can be used for immense good and therefore should be avoided, however, the negative impact it could have if not careful could disastrous and something I am very passionate about keeping a very close eye on.

### A.3.4 Professional issues

If this solution were to be used in a professional setting, such as recording chess games in competitions, and mistakes were made by the model, then I could see some pretty bad issues. I can't image Carlsen being very happy if he wanted to analyse his game but the inference application failed to record to PGN 100% correctly. This was not the aim of the project and more measures and effort would need to be taken to avoid consequences in professional settings.

# **Appendix B**

## **External Material**

<This appendix should provide a brief record of materials used in the solution that are not the student's own work. Such materials might be pieces of codes made available from a research group/company or from the internet, datasets prepared by external users or any preliminary materials/drafts/notes provided by a supervisor. It should be clear what was used as ready-made components and what was developed as part of the project. This appendix should be included even if no external materials were used, in which case a statement to that effect is all that is required.>