# Aligning Superhuman AI with Human Behavior: Chess as a Model System

Reid McIlroy-Young[1]    Siddhartha Sen[2]
Jon Kleinberg[3]    Ashton Anderson[1]

[1]Department of Computer Science, University of Toronto

[2]Microsoft Research

[3]Department of Computer Science, Cornell University

July 2020

## Motivation

Superhuman AI systems are increasingly prevalent

### Some domains humans are moving away from

- facial recognition, path finding, identifying photos of dogs

### But other areas will continue see human participation

- poker, chess, some business decisions

# Learning from Superhuman AI

In the domains where humans have been superseded but continue to participate, this raises the possibility that we could learn from them

# Learning from Superhuman AI

In the domains where humans have been superseded but continue to participate, this raises the possibility that we could learn from them, or collaborate with them.

## Superhuman Reasoning

Superhuman AI systems can be difficult for humans to understand, making it difficult for us to interact with or learn from them

## Research Question

How can we bridge the gap between the AI's behavior and ours?

## Requirements

1. Superhuman AI
2. Observed Humans
3. Diverse Humans

# Requirements

1. Superhuman AI
2. Observed Humans
3. Diverse Humans
4. Parametrized Humans

## Question

Predict the next move a human, at a specific skill level, will make during a chess game

## Lichess

Lichess.org is a popular, free, open-source chess platform with over 1 billion games in its database

# ELO Rating

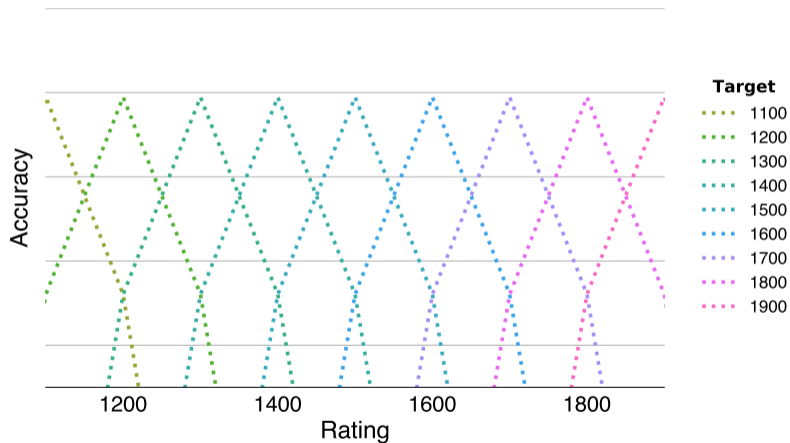ELO rating is a measurement of skill in a game, the larger the ELO rating the higher the skill

## Data

For testing we just use games from December 2019

1. Create bins for each range of 100 rating points
2. Divide games into the bins by the ELO of both players
3. Select 10,000 games from each bin, between 1000 and 2500

For each game we can then look at the mean move prediction accuracy of a model
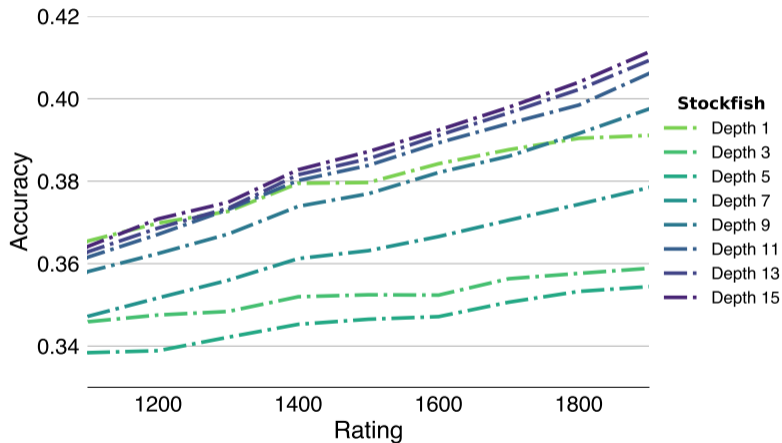
# What We Want

# Stockfish Overview

Design  Traditional Chess engine

Type  Tree Search

Humanity  Hand coded heuristics

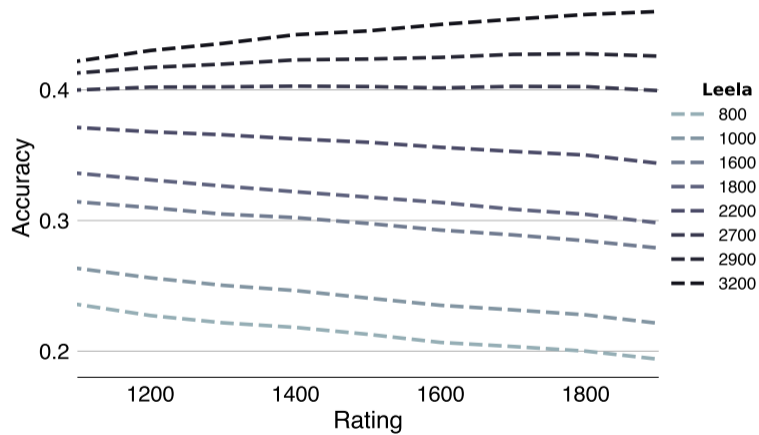Tuning  Depth of tree to search

# Stockfish

# Leela Intro

Design Implementation of AlphaZero

Type Reinforcement Learning

Humanity Only rules of chess
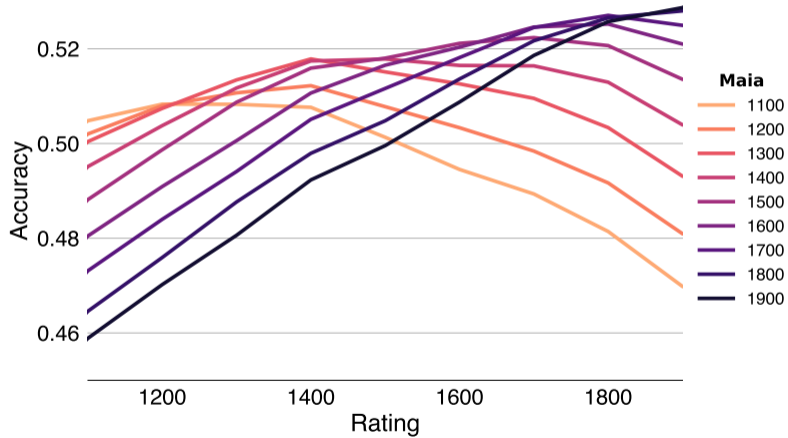
Tuning Length of training

# Leela

# Maia Intro

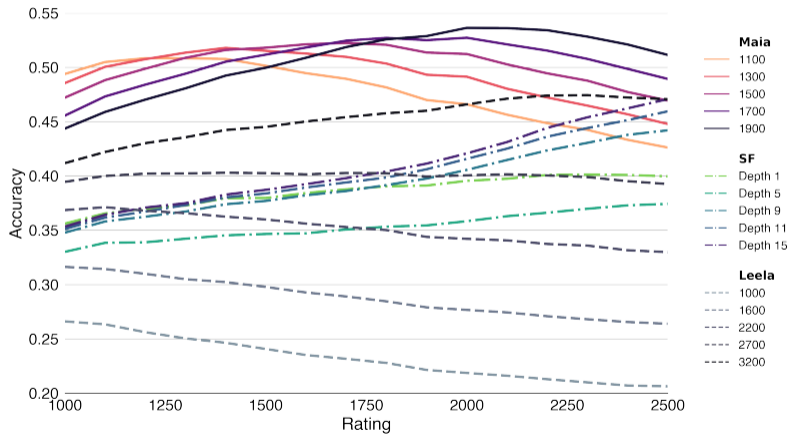Design  AlphaZero based deep neural net

Type  Classification

Humanity  Trained on 12 million human games each
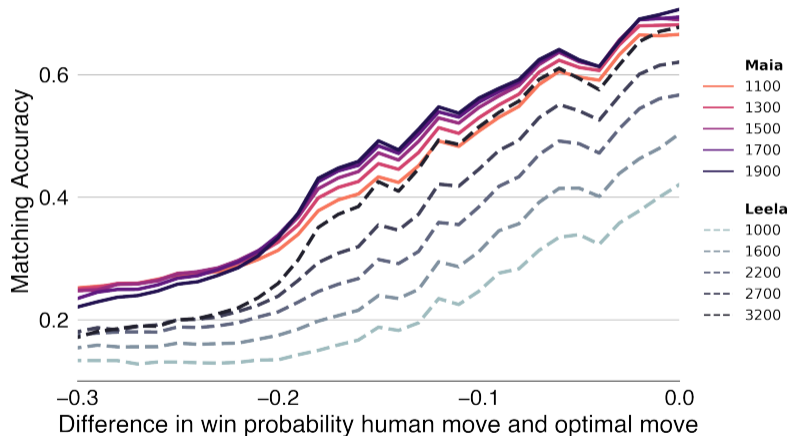
Tuning  Trained human skill

# Maia

# All

# Human Errors

## Discussion

Mistakes Categorizing

Understanding Human Skill

Learning Aid

## Further Information

Paper  KDD 2020

arXiv  arxiv.org/abs/2006.01855

Github  github.com/CSSLab/maia-chess

Lichess  maia1, maia5, maia9