# High Performance Data Compilation for Sports Analytics

Tyler Gourley
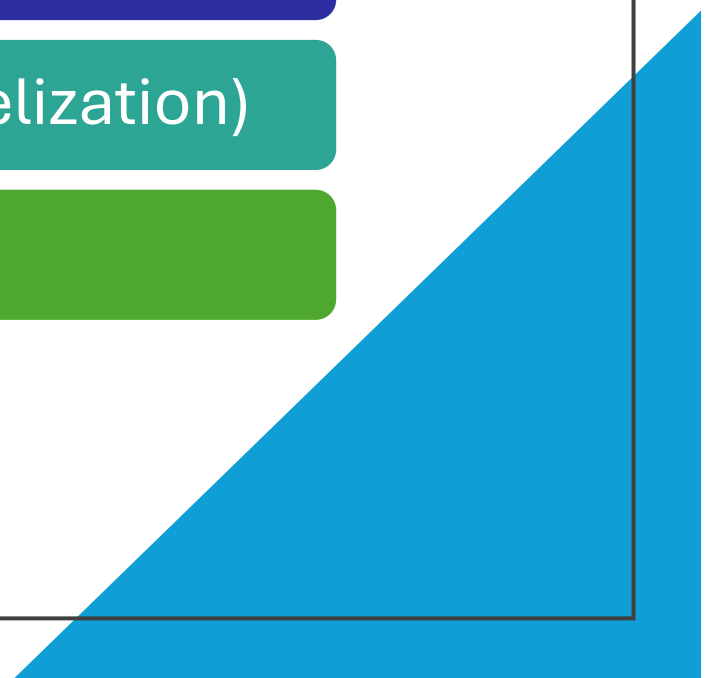
# Outline

- Background
- Optimization I (Set Intersection)
- Optimization II (Parallelization)
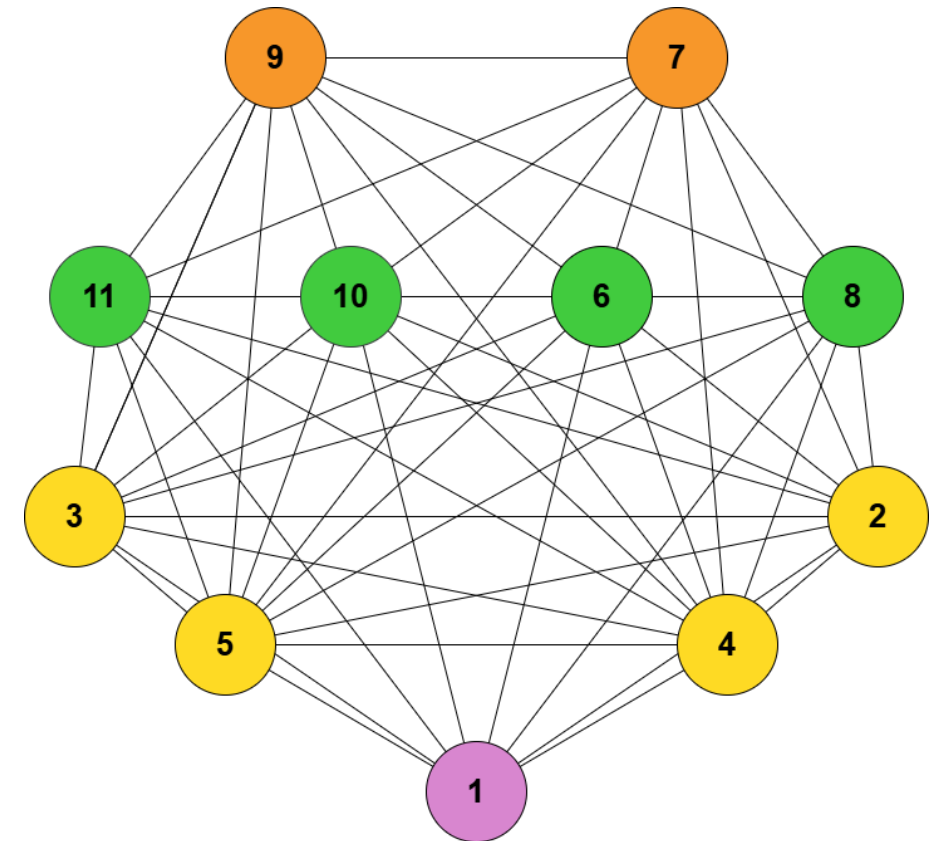- Results

# The Question...

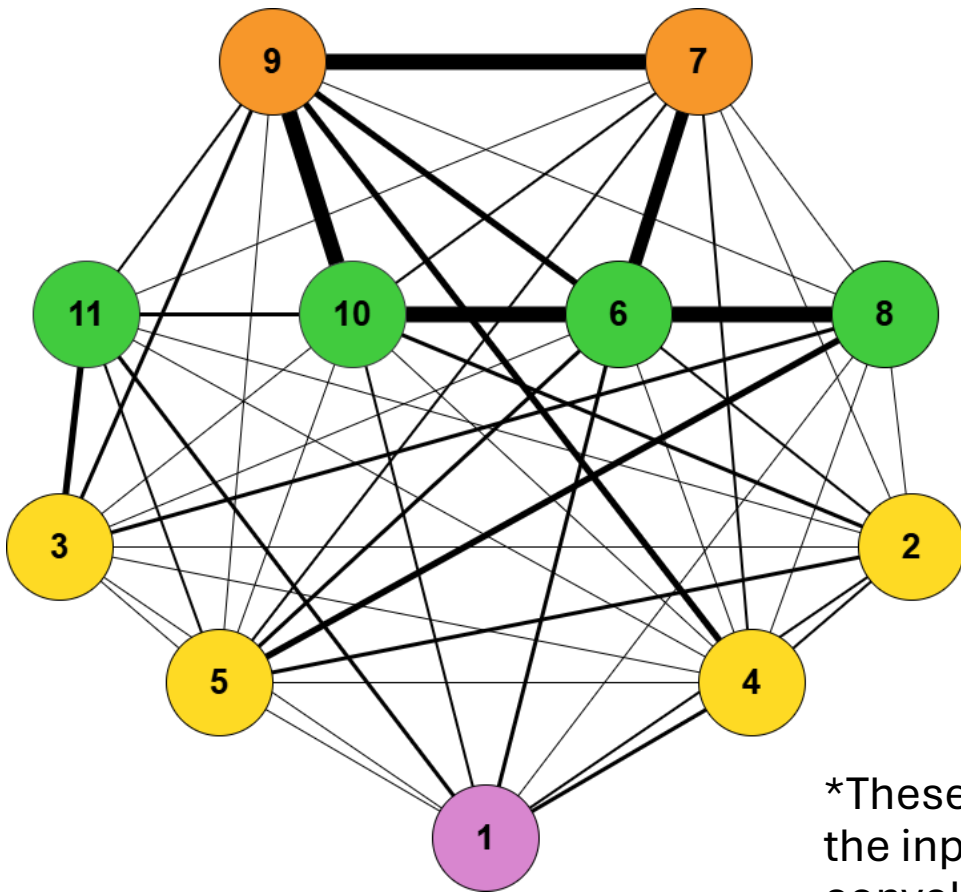Can we predict the outcome of soccer matches based on team chemistry?

# Graphical Representation

Each match will have 2 Complete, Weighted Graphs (one for each team)

# Numerical Representation

## Convert the Weighted Graphs to Numpy Arrays

*These arrays will be the input to our convolutional neural net

| 0 | 900 | 810 | 900 | 810 | 720 | 810 | 900 | 990 | 900 | 630 |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 0 | 0 | 810 | 810 | 810 | 720 | 630 | 810 | 720 | 540 | 360 |
| 0 | 0 | 0 | 720 | 630 | 990 | 900 | 720 | 810 | 810 | 630 |
| 0 | 0 | 0 | 0 | 810 | 990 | 900 | 810 | 1440 | 810 | 630 |
| 0 | 0 | 0 | 0 | 0 | 990 | 900 | 1530 | 990 | 810 | 630 |
| 0 | 0 | 0 | 0 | 0 | 0 | 2250 | 2160 | 1350 | 2250 | 630 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 900 | 2160 | 720 | 630 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 990 | 810 | 630 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2070 | 630 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 630 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

# Data Compilation Algorithm

**Algorithm 1:** Compile Team Chemistry Data

```
 1: matches ← query database (date_range)
 2: for match in matches
 3:     for team in match
 4:         init team_chem_array
 5:         for col in team_chem_array
 6:             for row in team_chem_array
 7:                 lineup_data1 ← query database (player_id1)
 8:                 for p1_data in lineup_data1
 9:                     lineup_data2 ← query database (player_id2)
10:                     for p2_data in lineup_data2
11:                         if p1_data.match_id = p2_data.match_id
12:                             team_chem_array[row][col] += 90
13:                         end if
14:                     end for
15:                 end for
16:             end for
17:         end for
18:         save team_chem_array
19:     end for
20: end for
```

loop iterations:

(330,000+)
(2)

(11)
(11)

(?)

(?)

Our database contains the playing time data for over 330,000 soccer matches

We iterate over every cell in the upper triangular array (representing the 55 player combinations possible)

For every cell, we need to query the database again to determine how many game minutes the players played together

# Data Compilation Algorithm

**Algorithm 1:** Compile Team Chemistry Data

loop iterations:

```
1:  matches ← query database (date_range)
2:  for match in matches                                (330,000+)
3:      for team in match                               (2)
4:          init team_chem_array
5:          for col in team_chem_array                  (11)
6:              for row in team_chem_array              (11)
7:                  lineup_data1 ← query database (player_id1)
8:                  for p1_data in lineup_data1          (?)
9:                      lineup_data2 ← query database (player_id2)
10:                     for p2_data in lineup_data2      (?)
11:                         if p1_data.match_id = p2_data.match_id
12:                             team_chem_array[row][col] += 90
13:                         end if
14:                     end for
15:                 end for
16:             end for
17:         end for
18:         save team_chem_array
19:     end for
20: end for
```

The bulk of the complexity comes from this inner 2 for loops

The results from these queries can be very large, and provide redundant information

We can do better

# Algorithm Optimization Using Sets

**Algorithm 1:** Compile Team Chemistry Data

loop iterations:

```
 1: matches ← query database (date_range)
 2: for match in matches                                    (330,000+)
 3:     for team in match                                   (2)
 4:         init team_chem_array
 5:         for col in team_chem_array                      (11)
 6:             for row in team_chem_array                  (11)
 7:                 lineup_data1 ← query database (player_id1)
 8:                 for p1_data in lineup_data1             (?)
 9:                     lineup_data2 ← query database (player_id2)
10:                     for p2_data in lineup_data2         (?)
11:                         if p1_data.match_id = p2_data.match_id
12:                             team_chem_array[row][col] += 90
13:                         end if
14:                     end for
15:                 end for
16:             end for
17:         end for
18:         save team_chem_array
19:     end for
20: end for
```

**Algorithm 2:** Compile Team Chemistry Data (Using Sets)

loop iterations:

```
 1: matches ← query database (date_range)
 2: for match in matches                                    (330,000+)
 3:     for team in match                                   (2)
 4:         init team_chem_array
 5:         L ← 0
 6:         for player_id in team                           (11)
 7:             S ← query database (player_id)
 8:             L[player_id] ← S
 9:         end for
10:         for row in team_chem_array                      (11)
11:             for col in team_chem_array                  (11)
12:                 R ← L[row] ∩ L[col]
13:                 team_chem_array[row][col] += 90 * |R|
14:             end for
15:         end for
16:         save team_chem_array
17:     end for
18: end for
19:
20:
```
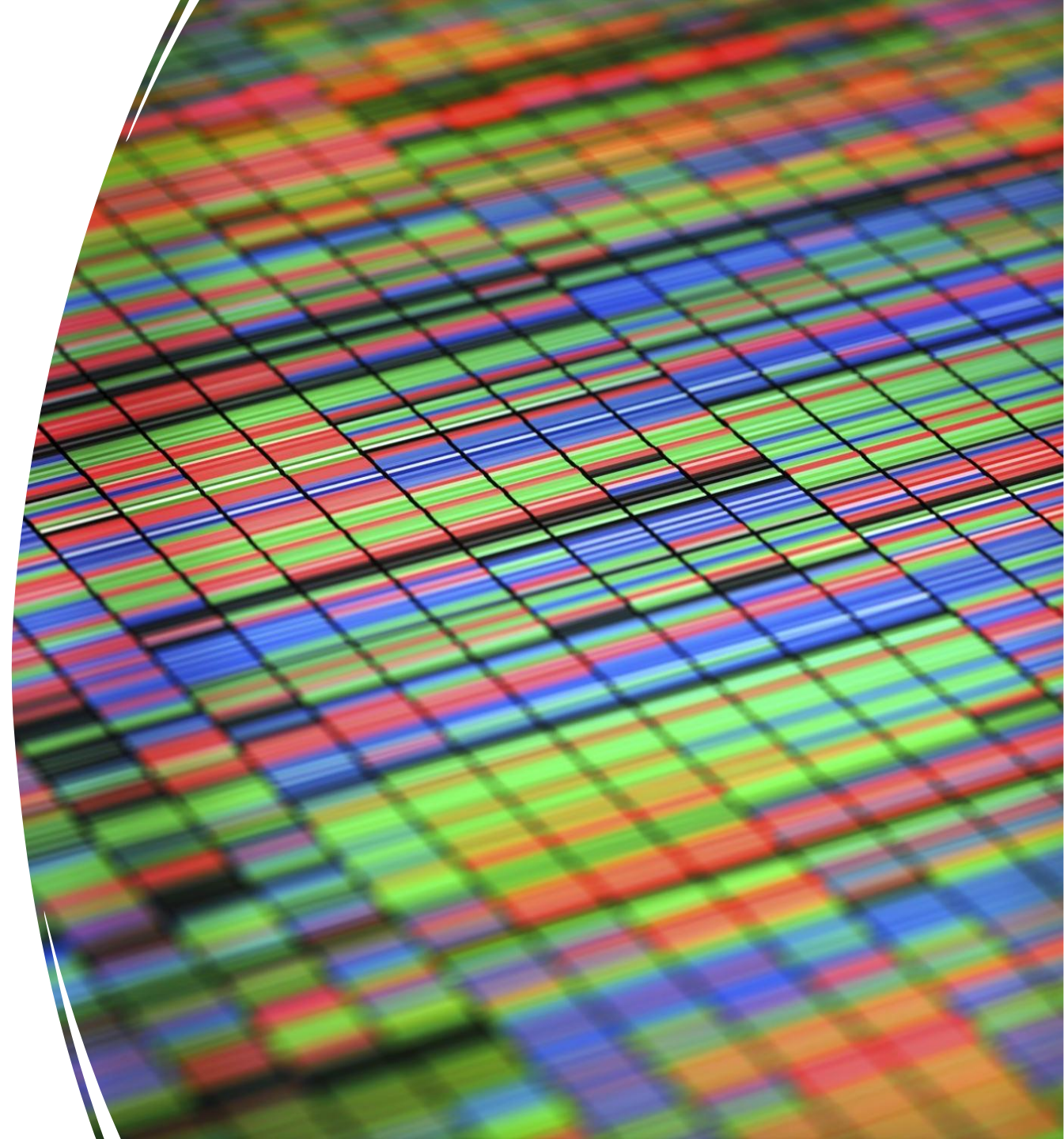
# Parallelization

Python Processes
(multiprocessing package)

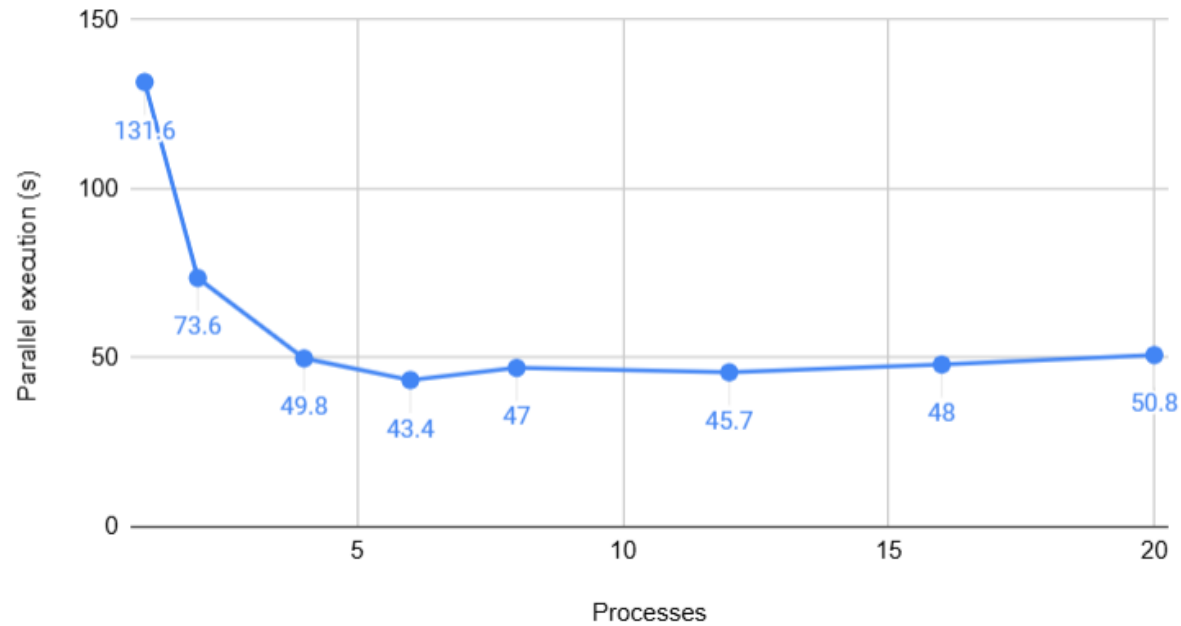Chunk the data into Tasks

Create a Task Pool

Add tasks to the pool

# Parallel Implementation Code

```python
# Query matches data from the database
matches = select_matches(start_date='2015-01-01', end_date='2024-01-01', limit=NUM_MATCHES)

# put the matches data (query result) into a pandas dataframe
df = pd.DataFrame(matches.fetchall())
df.columns = ["match_id", "date", "home_team_id", "away_team_id", "home_team_goal", "away_team_goal"]

# Divide the input data into chunks to send to the processes
chunk_size = 1 # if using a value other than 1, we need to modify our process_chunk function
chunks = [df.iloc[i:i + chunk_size] for i in range(0, df.shape[0], chunk_size)]
# print (chunks) # Uncomment to verify what the chunks look like

print(f'beginning parallel compilation')
print(f'compiling data for {NUM_MATCHES} matches')
print(f'with {NUM_PROCS} processors')

# Begin Parallelism
pool = multiprocessing.Pool(processes = NUM_PROCS) # create task pool
results = pool.imap_unordered(process_chunk, chunks)
pool.close() # signify that we are not adding any more tasks to the pool
pool.join() # blocking, waits for the entire task pool to be dried up
# End Parallelism
```
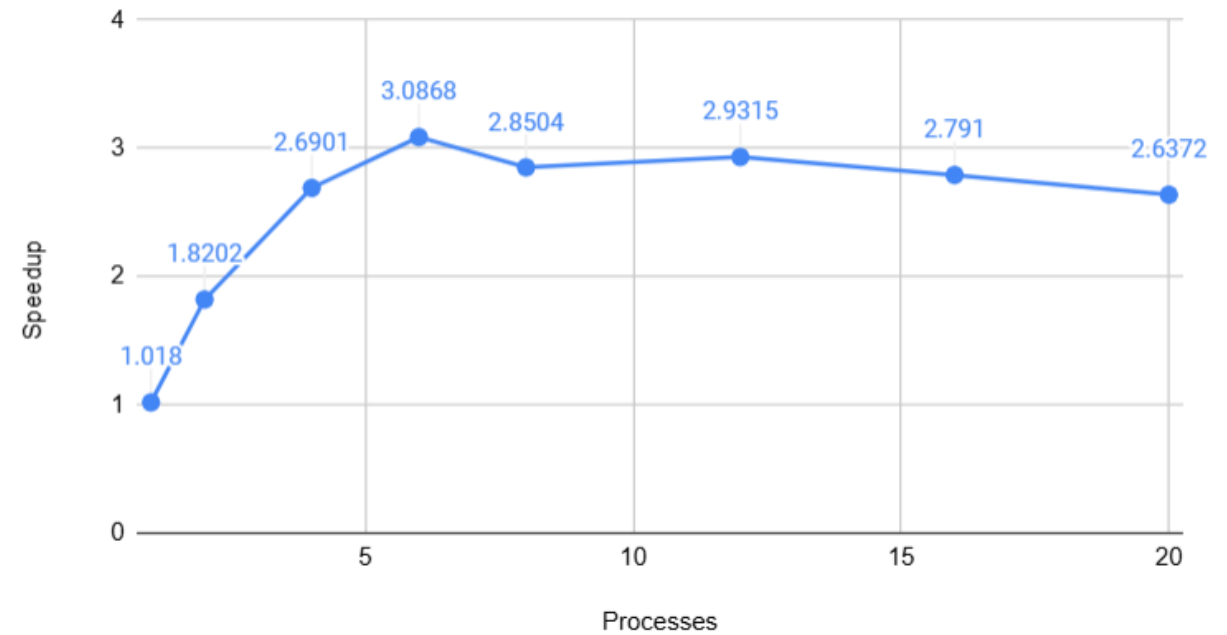
# Preliminary Results

The performance metrics were run on a representative 200 match portion of the overall database (330,000+ matches)
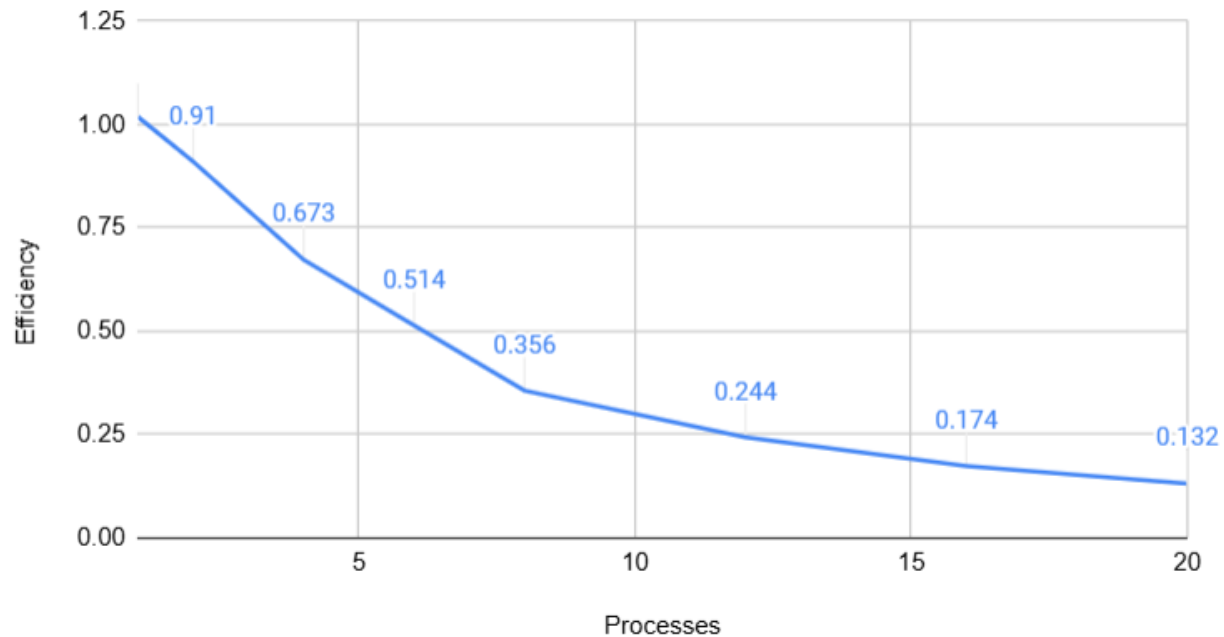


Team Chemistry - Parallel execution (s) vs. Processes



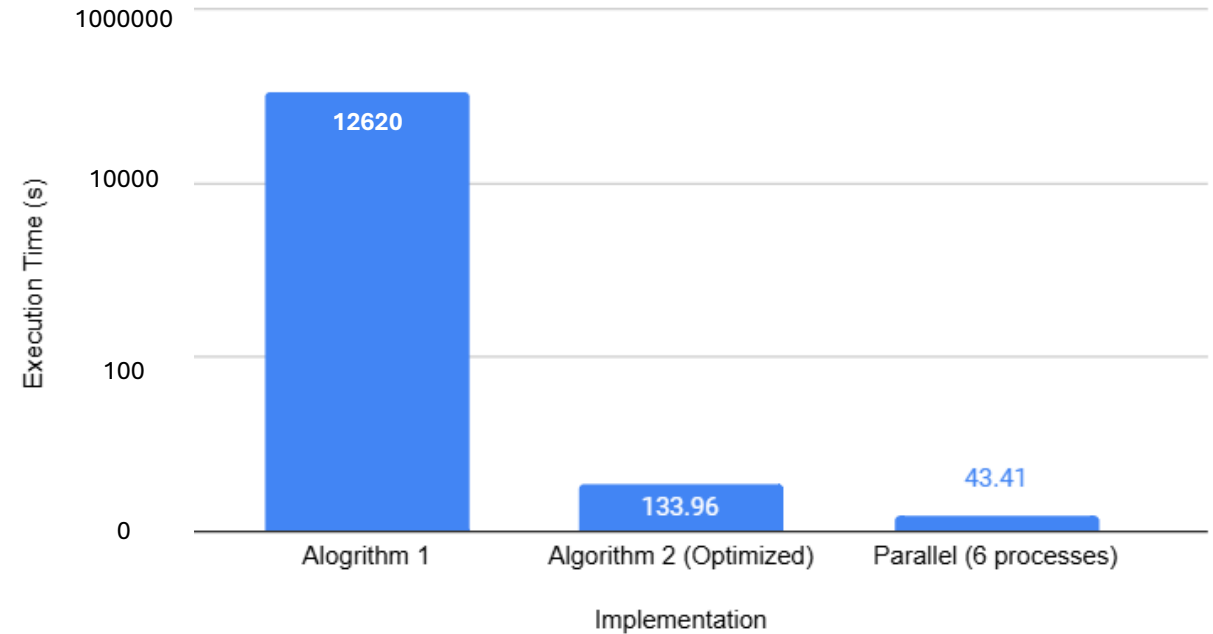Team Chemistry - Speedup vs. Processes

# Preliminary Results (Cont.)

We achieved a max Speedup of 3.086 when using 6 processes (about 300x faster when compared to original code)

## Team Chemistry - Efficiency vs. Processes



## Execution Time (s) vs. Implementation

# Questions?