# CS6023 Project - Matching with Preferences

J. Prem Krishnaa
CS14B049

R. Meenakshi
AE15B051

November 2018

## 1   The Stable Matching Problem

The stable matching problem has been very well studied in literature since its inception in the seminal paper by Gale and Shapley [1] in 1962. The stable matching problem has a lot of variants and researchers keep finding more and more pristine structural results on it. On the practical application side, we are interested in the problem as it can model many real world processes like college applications, student elective allocation etc. The problem is also studied from many different domains, for eg., researchers from economics study it from a varied perspective, with their interest arising from the recent popularity of matching markets. In this section we formally define the stable matching problem.

A stable matching instance can be modelled as a bipartite graph $G(\mathcal{M} \cup \mathcal{W}, E)$, where $\mathcal{M}$ represents the set of men and $\mathcal{W}$ represents the set of women. Every $x \in \mathcal{M} \cup \mathcal{W}$ has a strict preference order over a subset of elements in the other set. An edge $(m, w)$ is said to be acceptable, if $w$ is present in $m$'s preference list and vice-versa. The set of acceptable edges form the edge set $E$. Let $\mathcal{W}_m$ denote the set of women in $m$'s preference list and $\mathcal{M}_w$ denote the set of men in $w$'s preference list. We denote the strict preference order for $m$ by $>_w$ on $\mathcal{W}_m$, and similarly for $w$ by $>_w$ on $\mathcal{M}_w$.

Any arbitrary subset $M$ of $E$ is called an **matching**. For a given matching $M$, we denote the woman matched to $m$ by $M(m)$, and the man matched to $w$ by $M(w)$.

**Definition 1.** An edge $(m, w) \in E \setminus M$ is called a **blocking** pair if both the following conditions are satisfied: (i) $m$ is unmatched or $w >_r M(m)$, and (ii) $w$ is unmatched or $m >_w M(w)$ A matching is said to be **stable** if it does not admit any blocking pair.

For any stable matching instance, the Gale and Shapley algorithm, as described below, always computes a stable matching in $O(|E|)$ linear time. A lot of algorithms for different variants of stable matchings draw inspiration from Gale and Shapley algorithm. So if we are able to parallelize Gale and Shapley algorithm on GPU, it is very likely that we could be able to parallelize other similar algorithms with appropriate modifications. To make our task simpler, we make the following assumptions: the graph is a complete bipartite graph, i.e., all men have all women in their preference list and vice versa.

**Example 1.1.**
$$m_1 : w_1, w_2 \qquad\qquad w_1 : m_1, m_2$$
$$m_2 : w_2, w_1 \qquad\qquad w_2 : m_2, m_1$$

In the above example, we can verify that $M_1 = \{(m_1, w_1), (m_2, w_2)\}$ is a stable matching by our definition, whereas $M_2 = \{(m_1, w_2), (m_2, w_1)\}$ is not a stable matching.

**Parallelization:** A very interesting observation about Algorithm 1 is that in step-4, the order in which men make proposals is not specified. But irrespective of this order, the algorithm outputs the same final matching $M$. This property can be exploited for parallelization. So we can have men make proposals in parallel to their most preferred women, but a sequential bottleneck could occur, if all men apply to the same woman, in which case we need to process each man one by one using queue like structure. In terms of GPU implementation, we can consider it synonymous to an atomic operation.

**Algorithm 1** Gale and Shapley algorithm [1]

---

1: Input : $G = (\mathcal{M} \cup \mathcal{W}, E)$
2: Start with all $m \in \mathcal{M}$, $w \in \mathcal{W}$ as free, $M = \phi$
3: **do**
4:    **if** $\exists m \in \mathcal{M}, M(m) = \phi$, $\exists w \in \mathcal{W}_m$, s.t. $w$ is most preferred & not applied by $m$ **then**
5:       **if** $h$ is unmatched **then**
6:          $M = M \cup \{(m, w)\}$
7:       **else**
8:          Let $m'$ be the current partner of $w$
9:          **if** $m' >_w m$ **then**
10:             $m$ continues applying to next woman in his list (if it exists)
11:          **else**
12:             $M = (M \cup \{(m, w)\}) \setminus \{(m', w)\}$
13:             $m'$ continues applying to next woman in his list (if it exists)
14:    **else**
15:       Exit the loop.
16: **while** true
17: Return $M$

---

# 2 Serial code bench-marking

We had a public C++ repository (GraphMatching - Amit Rawat) containing stable matching serial code available to us. However since we felt it would be harder to port the existing code with CUDA, we implemented our own serial code using C language. We generated stable matching instances using three models as follows. These generators were described in Amit's python implementation repository (HR). Notation: $k$ is the fixed number of women in each man's preference list, for which in complete bipartite graph, $k = n$.

- **Master:** In this model we try to reflect a real world scenario that few women are more preferred among the men. Also the women have a master list to rank men. A geometric probabilistic distribution with $p = 0.10$ is set up among the women - this denotes the probability with which a woman is chosen to be included in a man's preference list. Each man samples $k$ women from the above distribution and orders them arbitrarily. Every woman orders the men in his preference list according to the master list of men.

- **Shuffle:** This model is similar to Master model except that we do not consider the master list of men. The men chose their preference list as described in Master model. But every woman orders the men in her preference list arbitrarily. This model also reflects a real-world scenario that each woman can have her own precedent for assigning preferences to men.

- **Random:** This model is used to generate arbitrary/unstructured instances by randomly sampling the preference list for men. Similarly every woman orders the men in her preference list arbitrarily.

In each model we considered the following values of $n$ (number of men, number of women): 10, 50, 100, 200, 500, 1000. Note that we are under a stable matching setting with complete bipartite graph / preference lists assumption. For each configuration, we measured the computation time and averaged it over five instances. As the algorithm has a running time of $O(|E|)$ and $|E| = n^2$, we expect to see a quadratic plot for $n$ vs computation time. The experimental time taken is plotted as follows.
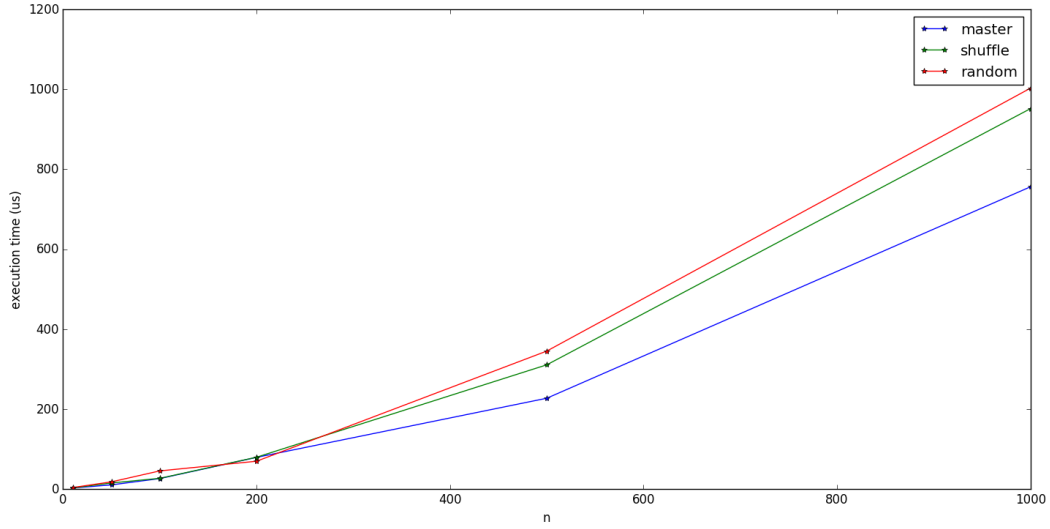
Figure 1: Serial - Time taken vs $n$

The above time taken are obtained by running on a machine with the following configuration: 4 GB RAM, Ubuntu 16.04, gcc compiler, Intel i5 processor. It is practically difficult to obtain the time taken for larger values of $n$. For eg., the size of one input file for $n = 10^4$, would be 1 GB. For $n = 10^5$, it would be 100 GB, which is impossible to process normally. It might be possible to obtain the time taken for $n = 10^4$, bu running on some cloud based service like AWS or Google Collab, since the current configuration is incapable of processing that.

Though we do not have too many data points, its easy to observe that the slope keeps increasing, which is similar to that of a quadratic curve. We can also see that there is not much difference between the execution times of different models. Master model is relatively faster, which could be attributed to the fact that it utilizes a master list of residents.

## 3 CUDA Implementation Attempts

### 3.1 CUDA Kernel-1

We implemented a basic CUDA kernel to compute a stable matching in a stable matching instance and verified its correctness by comparing with the output of serial implementation. We computed the time taken to run locally on a Nvidia GeForce 830M GPU card.

The idea behind parallel implementation was to have a separate thread for each man, which can apply to women independently of other men. For our preliminary purpose, one kernel call would perform exactly one proposal per unmatched man. We repeatedly perform the kernel call, until all the men are matched (this is guaranteed since we have equal number of men and woman and a complete bipartite setting.

As explained earlier in the parallelization intuition, we need to serialize the proposals happening to the same woman in a given time-step. For this purpose, we also implemented a lock using `atomicCAS` function. This could lead to performance degradation, as can be seen in the following plots - parallel implementation is almost 30 times slower. We are trying to see if we can perform some privatization, so that atomic operations are performed on shared memory instead of global memory. Also since our data set has a maximum value of $n = 1000$, we have currently used a kernel with 1 block and $n$ threads. We can try to tune these parameters as well to improve performance.
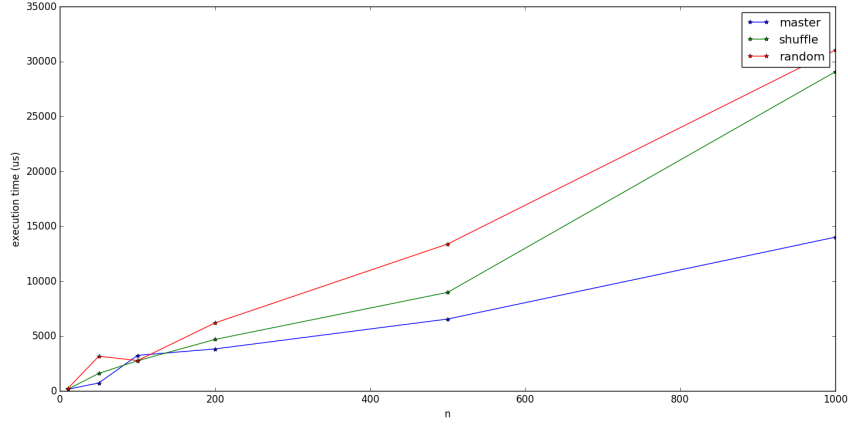
Figure 2: Kernel-1 : Time taken vs n

We also used Nvidia's profiling and visualization tools, `nvprof`, `nvvp` to check the critical time taking components of the CUDA implementation.
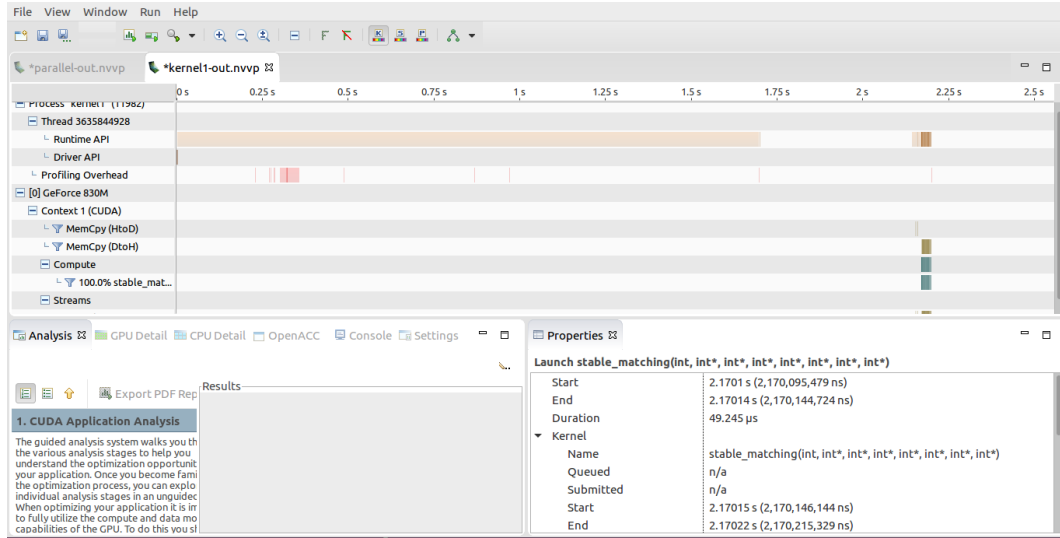


Figure 3: Kernel-1 Profiling

The above metrics were obtained for an instance with $n = 10^3$. The total compute time given by `cudaEvents` was 28 ms. However on checking the visual profiler, we found that the kernel was executing for only 10 ms. And the overhead due to `cudaMemcpy` from device to host in every iteration accounted for 1.5 ms. It turns out that the major part of the overhead is due to kernel launches (as can be seen in Figure 3).

## 3.2 CUDA Kernel-2

We suspected a bottleneck in Kernel-1 due to atomic add operations that we use to maintain the size of the matching, which could lead to inherent serialization. So we tried to remove this add operation and replace it with a boolean variable instead (non-atomic operation, to keep track of end of the algorithm). However this variable has to be set and checked before and after kernel call respectively, in every iteration, thus requiring

4

two `cudaMemcpy` operations, as compared to just one in Kernel-1. We verified the correctness of this kernel by comparing its output with that of serial code. Though we expected this kernel to perform better as we removed the atomic add operation, but in reality this was two times slower than Kernel-1, implying that the `cudaMemcpy` operations are much costlier.
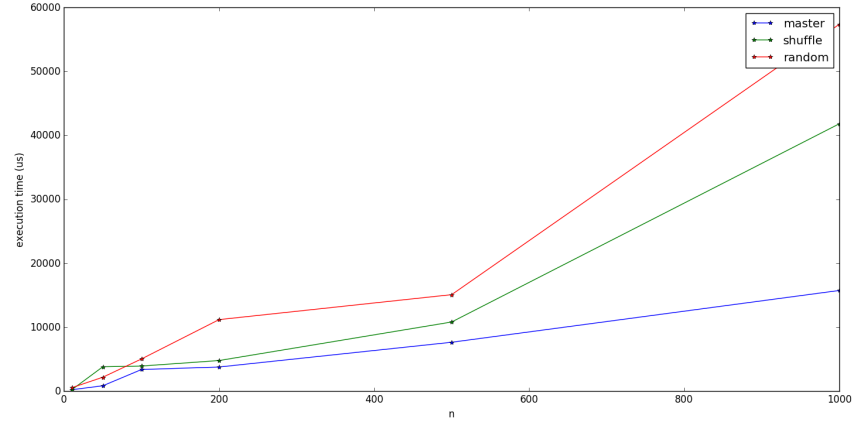


Figure 4: Kernel-2 : Time taken vs n

We used profiling tools `nvprof`, `nvvp` to check different components of our CUDA implementation.
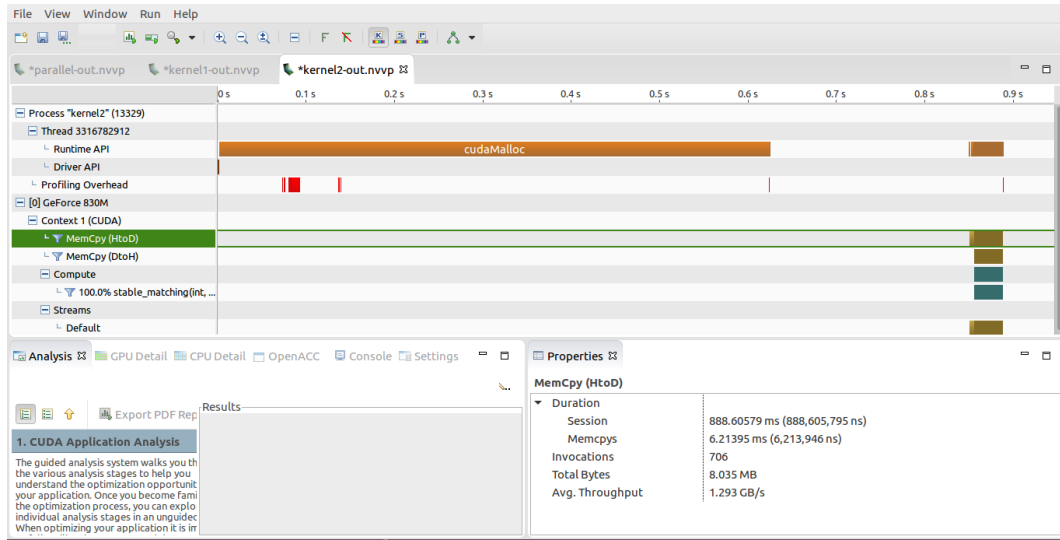


Figure 5: Kernel-2 Profiling

Turns out, atomic add overhead is negligible, and we are only adding more overheard due to extra `cudaMemcpy`. Further we still have the multiple kernel invocation overhead, carrying over from Kernel-1.

## 3.3   CUDA Kernel-3

We wanted to overcome the limitation of Kernel-2, by avoiding `cudaMemcpy` operations, which led us to think that if the main kernel call was done from the GPU in another kernel, then we need not have to do the

`cudaMemcpy` operations since the global variable is available in both the kernels. Hence we decided to try out **dynamic parallelism** by calling a `__global__` function from another `__global__` driver function. The driver function is called from the main with just one block and one thread configuration. Here we expected the runtime to be better than that of both Kernel-1 and Kernel-2, but however the runtime was almost similar to that of Kernel-1. The reason we believe is because of `cudaDeviceSynchronize` function needed to be placed in the driver kernel in order to start the child kernel, which is leading to an additional overhead.
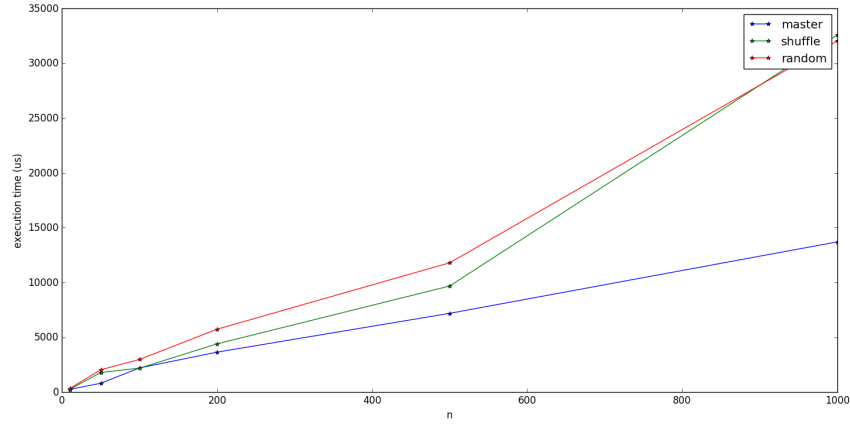


Figure 6: Kernel-3 : Time taken vs n

We used profiling tools `nvprof`, `nvvp` to check different components of our CUDA implementation.
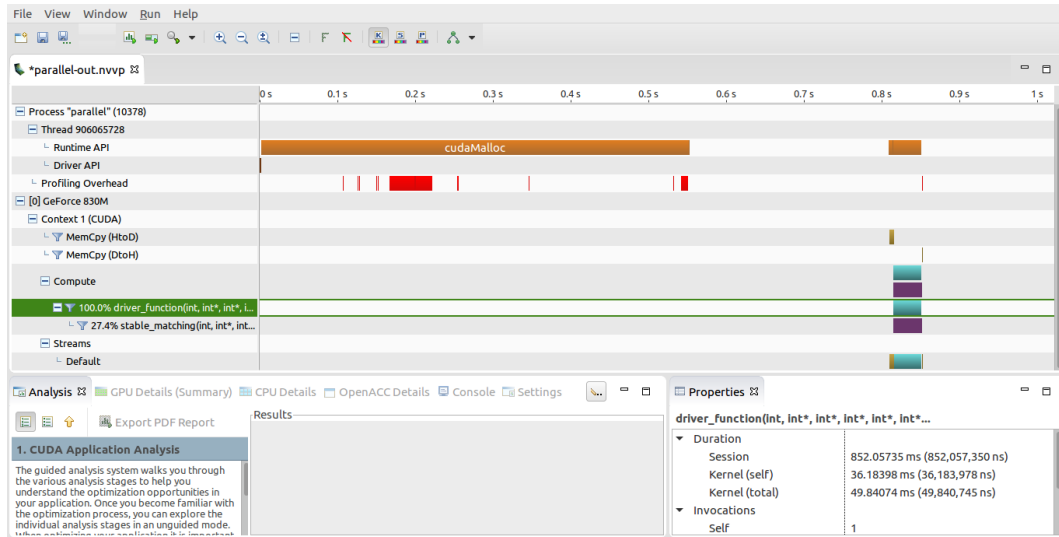


Figure 7: Kernel-3 Profiling

As we can observe from Figure 7, the actual kernel is taking up only 27% of the total time taken by the driver kernel. And the major part of the driver kernel is also taken up by `cudaEventSynchronize` command, as we predicted. This seems to be a drawback of using dynamic parallelism for our task.

6

# 4 Summary of Observations

Firstly we were limited by the maximum input size ($n = 10^3$) we could take due to compute limitation, as explained earlier. Our main focus in CUDA implementation was to speed up the compute time as compared to that of the serial implementation. Hence we did not take the reading time taken for serial into account. Also as far as we observed, there is no scope for pipe-lining the `cudaMemcpy` and kernel computation operations, since the entire input data is required before starting any kernel. Therefore asynchronous memory copies using streams would be of no use here. Also the time taken for memory copy is quite negligible compared to that of computation, at least for our input instances.

There is no global memory coalescing possible either, since the preference lists of men and woman are randomly generated. Inherently there are also lot of unavoidable `if` conditions in the algorithm, leading to control divergence between threads of a warp. We attempted to use privatization on top of our Kernel-1, but however as we observed from `nvvp` profiling, the major bottleneck is due to multiple kernel launches, and the overhead to initialize the shared memory for each launch would exceed the gain obtained over shared memory `atomicCAS` operation latency.

Finally in Kernel-3, we used dynamic parallelism, where we intuitively wanted to implement some sort of a master-slave model. But even there, we were unable to overcome the overhead due to `cudaDeviceSynchronize` calls, and ultimately leading to similar time taken as that of Kernel-1.

# 5 Conclusion

We tried to parallelize the Gale and Shapley algorithm [1] in CUDA, and we did so using three different kernel implementations. However despite our best efforts, the time taken for the parallel implementation turned out to be $\approx 30$ times slower than that of our most efficient serial implementation. Even though we believed at the start that our parallelization intuition could lead to an average case speed-up, it turned out not to be the case, due to the inherent overheads associated with the CUDA implementation such as kernel launch, atomic operations and `cudaDeviceSynchronize` calls in case of dynamic parallelism.

Hence with our attempts and experiments as evidence, we conclude that the Gale and Shapley algorithm is inherently **sequential**.

# References

[1] D. Gale and L. S. Shapley. College admissions and the stability of marriage. *The American Mathematical Monthly*, 69(1):9–15, 1962.