

Improvement Project : Parallel Gene Sequence Alignment

April 8, 2009 / CS 312 Algorithms

Duane Johnson

Objective

The Needleman-Wunsch algorithm as implemented in Project 5 is a serial solution to the gene sequence alignment problem. This paper reports on a parallel solution using Nvidia's CUDA extensions to C on a general-purpose graphics processing unit (GPGPU).

Solution

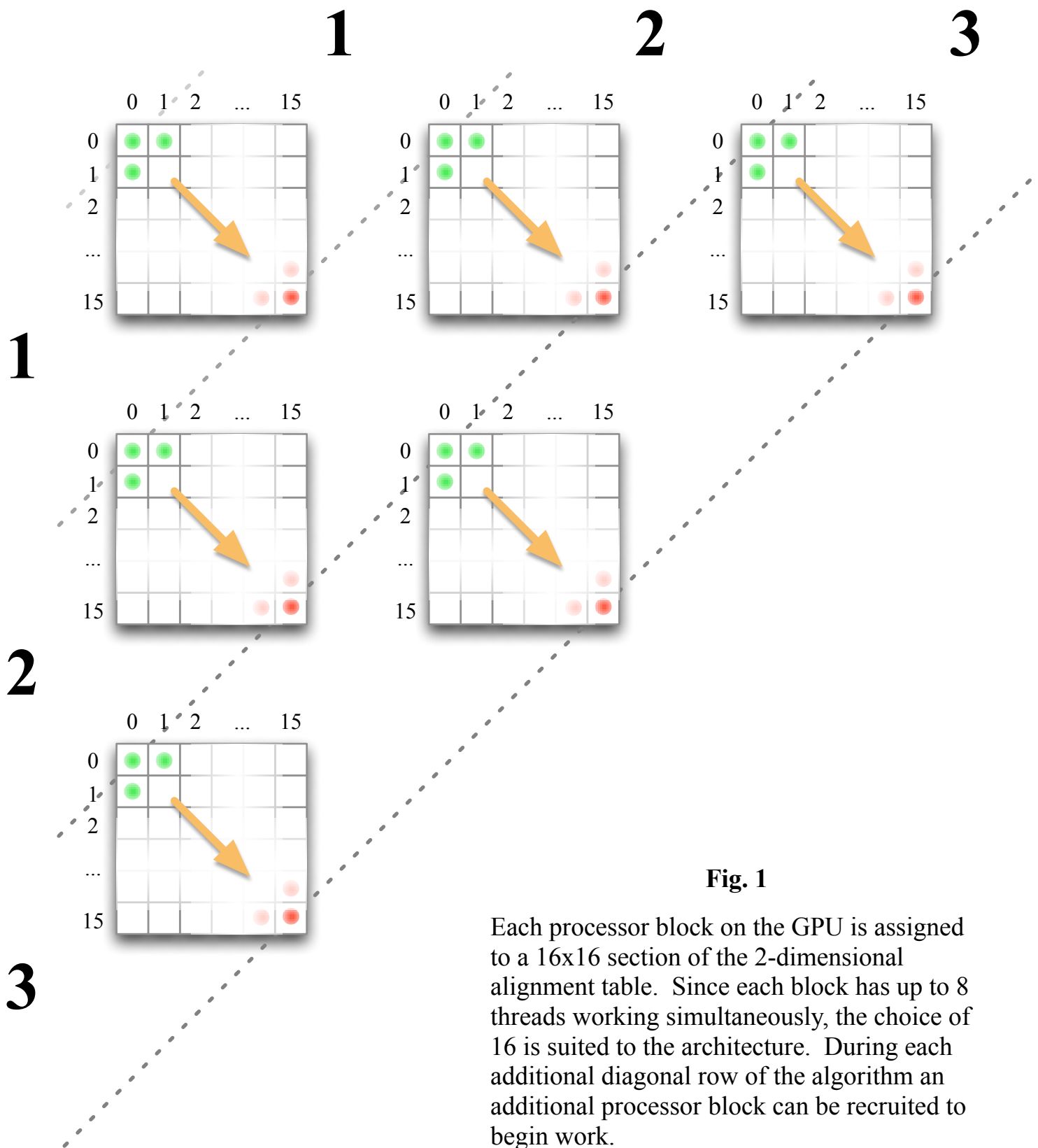
Instead of solving the problem in left-to-right and top-to-bottom fashion, this parallel algorithm works from the upper-left corner to the lower-right corner of the scoring table. It also partitions the problem into smaller units that are sized appropriately for the 240-core 8-series graphics processor from Nvidia. Once the scoring is achieved in parallel, the actual alignment backtrace proceeds as it did before--that is, in serial.

Details of the Parallel Algorithm

The 8-series graphics hardware consists of 30 blocks of processing cores, each block containing 8 cores (i.e. a total of 240 cores). In order to take advantage of these cores, it seemed reasonable to divide the scoring table up into sections of 16 x 16 and map these sections to blocks on the hardware (see Fig. 1).

The algorithm starts in the upper-left hand corner of the scoring table and assigns the first 16 x 16 section of the table to a block. The block's 8 cores then proceed in a diagonal upper-left-to-lower-right calculation step, adding up and scoring the various pathways as they go. Since each cell depends only on its three neighbors (one to the left, one above, and one diagonally to the left and above) the 8 cores are free to operate at near-optimal speed--"near optimal", that is, because they cannot all begin work at first, since each row must wait until the first cell on the row above it has been scored.

Once the first 16 x 16 section has been calculated, the section to its right and the section below it can begin to be calculated in parallel. Thus, the second iteration of 16 x 16 sections (i.e. a diagonal "row") can potentially be achieved as fast as the first iteration. This block-wise diagonal marching continues until the entire table is scored and an optimal score is entered in the lower-right hand cell of the table.



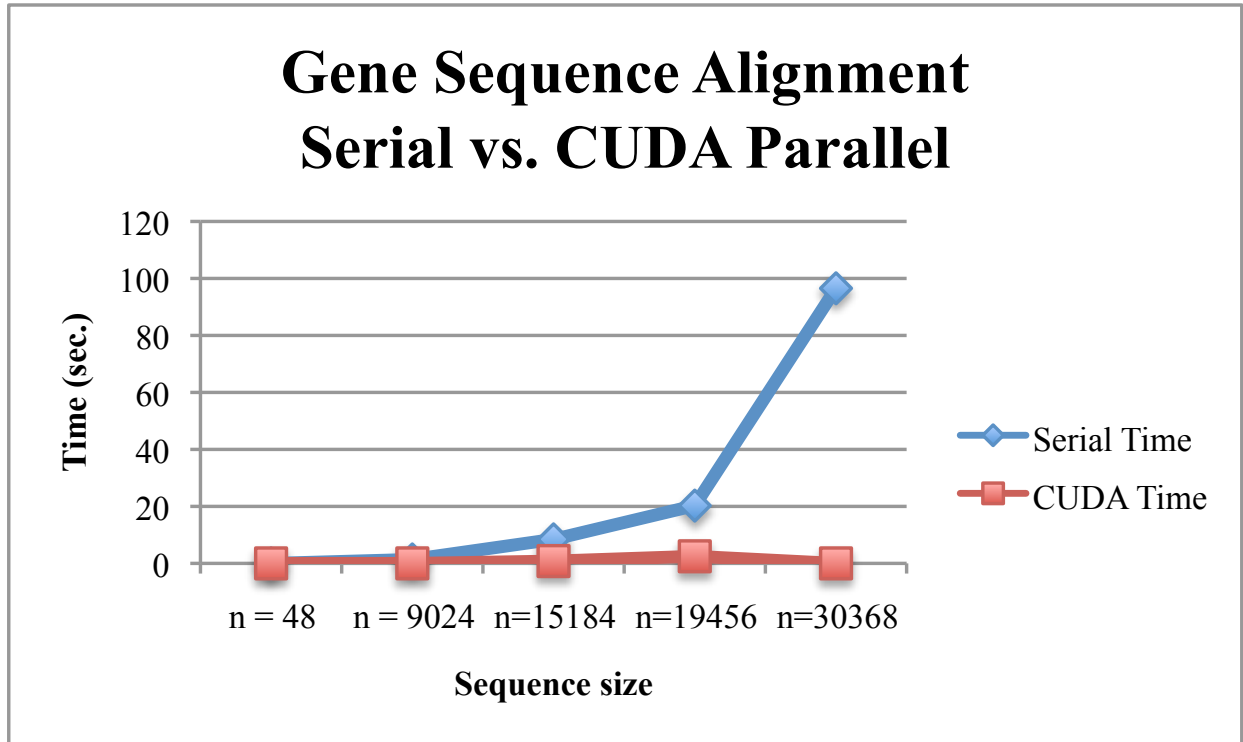


Fig. 2 Parallel alignment yields a speedup of about 10 X

	n = 48	n = 9024	n=15184	n=19456	n = 30368
Serial Time	0.000018	1.392797	8.435221	20.276045	96.541793
CUDA Time	0.000043	0.005645	0.813134	2.487627	failed
Space	2 kB	81 MB	231 MB	379 MB	922 MB

Table 1 Measurements of scoring times from n = 48 to n = 30368

Analysis

Ignoring setup time for the GPU, the absolute speedup is nearly 10 X for values of n in the range $15,000 < n < 20,000$. Part of this speedup is due to the fact that the serial implementation is not optimized for serial processing; nevertheless, the performance improvement is significant for problems whose scoring table can be fit in the GPU's memory. As can be seen in the case of n = 30368, the CUDA score failed since the table could not be transferred to GPU memory.

Since many alignment problems are in the range of thousand or tens of thousands of base pairs, parallel-accelerated speedups using GPGPUs is a practical way to reduce alignment times for sizes in this range by an order of ten.