

Green-Marl: A Domain-Specific Language for Easy and Efficient Graph Analysis

Oct 25th 2023
Wenxuan Li

Context & Motivation



- Context: large amount of graph data to be analysed and mined
- Challenges for efficient graph analysis on large-scale graph data:
 - Capacity: limited memory space
 - **Performance**: when run on large scale data
 - **Implementation**: difficult to implement correctly and efficiently
- Main focus: tight coupling between high-level graph analysis algorithm design and underlying hardware architecture

Overview



Green-Marl:

- A high-level domain specific **language**
- An associated **compiler** for producing optimized and parallelized low-level implementation

=> for both performance (optimization) and implementation (decoupling)

Language Design: Domain-specific Syntaxes

```
1 Procedure Compute_BC(  
2   G: Graph, BC: Node_Prop<Float>(G)) {  
3   G.BC = 0; // initialize BC  
4   Foreach(s: G.Nodes) {  
5     // define temporary properties  
6     Node_Prop<Float>(G) Sigma;  
7     Node_Prop<Float>(G) Delta;  
8     s.Sigma = 1; // Initialize Sigma for root  
9     // Traverse graph in BFS-order from s  
10    InBFS(v: G.Nodes From s) (v!=s) {  
11      // sum over BFS-parents  
12      v.Sigma = Sum(w: v.UpNbrs) {w.Sigma};  
13    }  
14    // Traverse graph in reverse BFS-order  
15    InRBFS(v!=s) {  
16      // sum over BFS-children  
17      v.Delta = Sum (w:v.DownNbrs) {  
18        v.Sigma / w.Sigma * (1+ w.Delta)  
19      };  
20      v.BC += v.Delta @s; //accumulate BC  
21    } } }
```

- Data-Types:
 - Graph
 - Nodes
 - Node_Prop
- Traversal:
 - InBFS (BFS)
 - InRBFS (reverse-order BFS)
 - UpNbrs and DownNbrs

Figure 1. Betweenness Centrality algorithm described in Green-Marl

Language Design: Parallelism

```
1 Procedure Compute_BC(  
2   G: Graph, BC: Node_Prop<Float>(G) {  
3     G.BC = 0; // initialize BC  
4     Foreach(s: G.Nodes) {  
5       // define temporary properties  
6       Node_Prop<Float>(G) Sigma;  
7       Node_Prop<Float>(G) Delta;  
8       s.Sigma = 1; // Initialize Sigma for root  
9       // Traverse graph in BFS-order from s  
10      InBFS(v: G.Nodes From s) (v!=s) {  
11        // sum over BFS-parents  
12        v.Sigma = Sum(w: v.UpNbrs) {w.Sigma};  
13      }  
14      // Traverse graph in reverse BFS-order  
15      InRBFS(v!=s) {  
16        // sum over BFS-children  
17        v.Delta = Sum (w:v.DownNbrs) {  
18          v.Sigma / w.Sigma * (1+ w.Delta)  
19        };  
20        v.BC += v.Delta @s; //accumulate BC  
21    } } }
```

- Implicit Parallelism
 - Group Assignment: `G.BC=0`
- (Explicit) Parallel Execution Region
 - `Foreach`
 - following fork-join style

Figure 1. Betweenness Centrality algorithm described in Green-Marl

Compiler

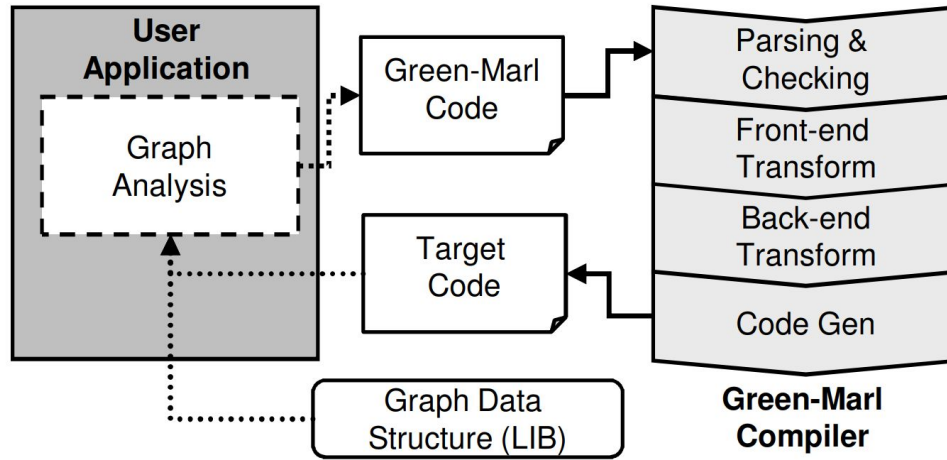


Figure 3. Overview of Green-Marl DSL-compiler Usage

Compiler: Loop Fusion



```
103  Foreach (s: G.Nodes) (f(s))
104      s.A = X(s.B);
105  Foreach (t: G.Nodes) (g(t))
106      t.B = Y(t.A)
```

becomes

```
107  Foreach (s: G.Nodes) (
108      if (f(s)) s.A = X(s.B);
109      if (g(s)) s.B = Y(s.A);
110  }
```

Compiler: Set-Graph Loop Fusion

```
139  Node_Set S (G); // ...
140  Foreach (s: S.Items)
141      s.A = x (s.B);
142  Foreach (t: G.Nodes) (g (t))
143      t.B = y (t.A)
```

becomes

```
144  Foreach (s: G.Nodes) (
145      if (S.Has (s)) s.A = x (s.B);
146      if (g (s)) s.B = y (s.A);
147  }
```


Compiler: Code Generation and Architecture Portability



Compiler emits out target code using code-generation templates

- Example: `Foreach` implementation with backend OpenMP

```
222 Foreach(s:G.Nodes)
223     For(t: s.Nbrs)
224         s.A = s.A + t.B;
```

becomes

```
225 OMP(parallel for)
226 for(index_t s = 0; s < G.numNodes(); s++) {
227     // iterate over node's edges
228     for(index_t t=G.edge_idx[s]:t<G.edge_idx[s+1];t++){
229         // get node from the edge
230         index_t t = G.node_idx[t];
231         A[s] = A[s] + B[t];
232     } }
```

Allows replacement of backends for code-generations, e.g. CUDA

Evaluation: Productivity



- Measured by Line of Codes (LoC)
- Compared with implementations in existing graph analysis libraries

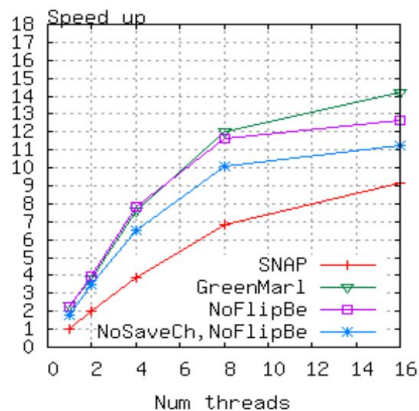
Name	LOC	LOC	Source
	Original	Green-Marl	
BC	350	24	[9] (C OpenMp)
Conductance	42	10	[9] (C OpenMp)
Vetex Cover	71	25	[9] (C OpenMp)
PageRank	58	15	[2] (C++, sequential)
SCC(Kosaraju)	80	15	[3] (Java, sequential)

Table 3. Graph algorithms used in the experiments and their Lines-of-Code(LOC) when implemented in Green-Marl and in a general purpose language.

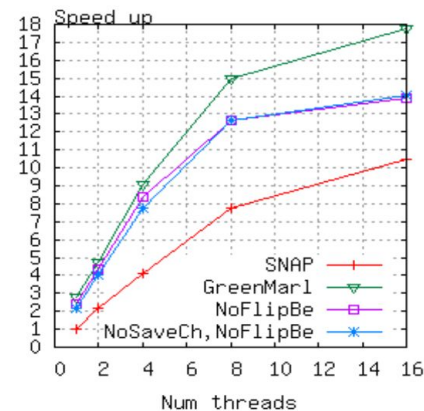
Evaluation: Performance Gain



- Measured by Speed-up with number of threads growing
- Compared with implementations in SNAP library (Bader et al. 2008)
- Ablation study by disabling some optimizations of the compiler (e.g. FlippingEdge)



(a) RMAT



(b) Uniform

Figure 4. Speed-up of Betweenness Centrality. Speed-up is over the SNAP library [9] version running on a single-thread. NoFlipBE and NoSaveCh means disabling the *Flipping Edges* (Section 3.3) and *Saving BFS Children* (Section 3.5) optimizations respectively.

(Evaluated on randomly synthesized graphs with 32 million nodes and 256 million edges)

Limitations



- No backends supported for distributed environments when the paper was released
 - Later works introduced Pregel (Hong et al. 2014), CUDA (Shashidhar and Nasre. 2017) and MPI (Rajendran and Nandivada 2020)
- No baseline provided for evaluating speed-ups on PageRank and Kosaraju's algorithm
- Still extra cost for mastering this language

Reference



- Bader et al. 2008. SNAP, Small-world Network Analysis and Partitioning: An open-source parallel graph framework for the exploration of large-scale networks
- Hong et al. 2014. Simplifying Scalable Graph Processing with a Domain-Specific Language
- Shashidhar and Nasre 2017. LightHouse: An Automatic Code Generator for Graph Algorithms on GPUs
- Rajendran and Nandivada 2020. DisGCo: A Compiler for Distributed Graph Analytics



Thank you