

15754 Project

Salil Joshi & Cyrus Omar

April 15, 2011

1 Grammar

We will work on a simple language that is essentially a simple subset of CUDA. Types:

$$\begin{aligned} T' &:= \text{uint} \mid \text{int} \mid \text{float} \mid \text{bool} \\ T &:= \text{global } T'[] \mid \text{local } T'[] \mid T' \mid \text{const } T' \mid \text{TID}(a, b) \end{aligned}$$

Terms

$$\begin{aligned} e &:= x \mid 1 \dots \mid \text{tid} \mid \text{bid} \mid a[e] \\ \text{stmt} &:= \langle \rangle \mid \text{cmd}; \text{stmt} \\ \text{cmd} &:= \text{if } e \text{ then stmt else stmt} \\ &\quad \mid \text{for } x \text{ in } [e, e, e] \{ \text{stmt} \} \\ &\quad \mid \text{sync} \\ &\quad \mid x := e \\ &\quad \mid a[e] := e \\ &\quad \mid \text{decl } T \ x \text{ in } S \end{aligned}$$

The language is a simple imperative language with a few special features:

- *tid* is special variable that holds the thread index of the current thread. Similarly, *bid* holds the block index.
- Arrays are either `global` or `local`. Global arrays are stored on the global device memory while local arrays are stored on the block level shared memory
- A block level synchronization primitive, **sync** which corresponds to CUDA's `__syncthreads__`
- `const` types. The value of an expression of this type can be computed at compile time.
- $\text{TID}(a, b)$. An expression of this type has the value $a * \text{tid} + b$ i.e. it is an affine transformation of the thread index

2 Typing

We have three typing judgements:

- $\Gamma/\Delta \vdash e : T$ for expressions e states that e has type T in the context Γ . The context Δ collects the set of indices read by e for every global array that e accesses. However, only those indices which are *affine transformations of the thread index* (i.e. have type $\text{TID}(a, b)$) are collected. Δ will later be used to determine the memory segments that are accessed by the program, which is useful for optimizations like memory coalescing and data sharing. At present we focus on memory coalescing.
- $\Gamma/\Delta_r/\Delta_w \vdash C \text{ ok}$ for commands C states that the command is well formed and all subexpressions are well typed. In addition, Δ_r is the set of global array indices *read from* by C and Δ_w is the set of global array indices *written to* by C , similar to Δ in the previous judgement. Once again we only track array indices that are affine transformation of the thread index.
- $\Gamma/\Delta_r/\Delta_w \vdash S \text{ ok}$ similar to the previous judgement, but for statements.

2.1 Expressions

$$\begin{array}{c}
\frac{x : T \in \Gamma}{\Gamma \vdash x : T} \qquad \frac{}{\Gamma \vdash \text{tid} : \text{TID}(1, 0)} \qquad \frac{}{\Gamma \vdash \text{bid} : \text{uint}} \\
\\
\frac{\Gamma/\Delta_r^1 \vdash e_1 : \text{TID}(a_1, b_1) \quad \Gamma/\Delta_r^2 \vdash e_2 : \text{TID}(a_2, b_2)}{\Gamma/\Delta_r^1, \Delta_r^2 \vdash e_1 + e_2 : \text{TID}(a_1 + a_2, b_1 + b_2)} \qquad \frac{\Gamma/\Delta_r^1 \vdash e_1 : \text{const uint} \quad \Gamma/\Delta_r^2 \vdash e_2 : \text{TID}(a, b)}{\Gamma/\Delta_r^1, \Delta_r^2 \vdash e_1 * e_2 : \text{TID}(a * e_1, b * e_1)} \\
\\
\frac{\Gamma \vdash a : \text{global } T'[] \quad \Gamma/\Delta_r \vdash e : \text{uint}}{\Gamma/\Delta_r \vdash a[e] : T'} \qquad \frac{\Gamma/\Delta_r \vdash a : \text{global } T'[] \quad \Gamma/\Delta_r \vdash e : \text{TID}(a, b)}{\Gamma/(a', (a, b)), \Delta_r \vdash a[e] : T'} \\
\\
\frac{\Gamma/\Delta \vdash e_1 : \text{uint}}{\Gamma/\Delta \vdash e_1 : \text{int}} \qquad \frac{\Gamma/\Delta \vdash e_1 : \text{const uint}}{\Gamma/\Delta \vdash e_1 : \text{TID}(0, e_1)} \qquad \frac{\Gamma/\Delta \vdash e_1 : \text{const } T'}{\Gamma/\Delta \vdash e_1 : T'}
\end{array}$$

2.2 Commands

$$\begin{array}{c}
\frac{x : T, \Gamma/\Delta_r/\Delta_w \vdash S \text{ ok}}{\Gamma/\Delta_r/\Delta_w \vdash \text{decl } T \text{ } x \text{ in } S \text{ ok}} \qquad \frac{x : T \in \Gamma \quad T_{\text{scalar}} \quad \Gamma/\Delta_r \vdash e : T}{\Gamma/\Delta_r \vdash x := e \text{ ok}} \qquad \frac{}{\text{sync ok}} \\
\\
\frac{\Gamma/\Delta_r^1 \vdash e : \text{bool} \quad \Gamma/\Delta_r^2/\Delta_w^2 \vdash S_1 \text{ ok} \quad \Gamma/\Delta_r^3/\Delta_w^3 \vdash S_2 \text{ ok}}{\Gamma/\bigcup \Delta_r / \bigcup \Delta_w \vdash \text{if } e \text{ then } S_1 \text{ else } S_2 \text{ ok}} \\
\\
\frac{x : T \in \Gamma \quad T = \text{global } T'[] \quad \Gamma/\Delta_r^1 \vdash e_1 : \text{TID}(a', b') \quad \Gamma/\Delta_r^2 \vdash e_2 : T' \quad \Gamma/\Delta_r^3/\Delta_w \vdash S \text{ ok}}{\Gamma/\bigcup \Delta_r / (a, (a', b')), \Delta_w \vdash a[e_1] := e_2 \text{ ok}} \\
\\
\frac{\Gamma/\delta_r^3 \vdash e_3 : \text{int} \quad \Gamma/\Delta_r^0/\Delta_w^0 \vdash S[e_1/x] \text{ ok} \quad \Gamma/\delta_r^1 \vdash e_1 : \text{int} \quad \Gamma/\delta_r^2 \vdash e_2 : \text{int} \quad \Gamma/\Delta_r^0/\Delta_w^0 \vdash S[e_1 + e_3/x] \text{ ok} \quad \dots \quad \Gamma/\Delta_r^{15}/\Delta_w^{15} \vdash S[e_1 + 15 * e_3/x] \text{ ok}}{\Gamma/\bigcup \Delta_r \cup \bigcup \delta_r / \bigcup \Delta_w \vdash \text{for } x \text{ in } [e_1, e_2, e_3]\{S\} \text{ ok}}
\end{array}$$

The rule for the **for** essentially unrolls the first 16 iterations of the loop. If the loop variable is being used in an array index this rule collects the first 16 values that the index can take. The reason is that the same behaviour repeats itself as far as memory coalescing is concerned after the first 16 iterations as the difference in addresses is a multiple of 16 (and there are 16 threads in a half warp)

2.3 Statements

$$\frac{\Gamma/\Delta_r/\Delta_w^1 \vdash C \text{ ok} \quad \Gamma/\Delta_r^2/\Delta_w^2 \vdash S \text{ ok}}{\Gamma/\Delta_r^1, \Delta_r^2/\Delta_w^1, \Delta_w^2 \vdash C; S \text{ ok}} \qquad \frac{}{\Gamma/\cdot/\cdot \vdash \langle \rangle \text{ ok}}$$