

## Team 1 Planning Document

Team 1(이진우, 정덕인, 김도현, 임지규)

### I. Pipeline for LLVM IR Optimization & Assembler

Blueprint of our optimizer pass pipeline is based on *clang -O1*. However, we removed some transform passes which are not necessary or has insignificant effect in our project specification while adding some passes which are expected to reduce the cost significantly. Custom passes we will introduce in our projects are written in red below. The basic information of each custom passes are introduced below, while details will be described in the 『Requirements and Specification』 document as they are implemented. Simple references for pre-existing passes could be found on [llvm.org/docs/Passes.html](http://llvm.org/docs/Passes.html).

1 ⇒ 2 ⇒ 3-1 ⇒ 4 ⇒ 5-1 ⇒ 4 ⇒ 5-2 ⇒ 4 ⇒ 6
---

- |  |
|--|
| <ol style="list-style-type: none"><li>1. GV Optimization &amp; Other global-level preprocessing<ul style="list-style-type: none"><li>- (InterProcedural) Sparse Conditional Constant Propagation + Simplify CFG</li><li>- GV optimizer</li><li>- Dead Global Elimination</li><li>- GV to malloc in main</li><li>- Mem2Reg</li><li>- Global Value Numbering</li></ul></li><li>2. Function Call Optimization<ul style="list-style-type: none"><li>- Function Merging</li><li>- Tail Call Elimination</li><li>- Function Inline</li></ul></li><li>3. Memory Operation Optimization<ul style="list-style-type: none"><li>- Vectorize Load/Store</li></ul></li><li>3-1. Stack/Heap Optimizations<ul style="list-style-type: none"><li>- Heap to Stack+Argument Lowering</li><li>- Interprocedural Alloca Sinking</li><li>- Stack/Heap Access Grouping</li><li>- Mem2Reg</li></ul></li><li>4. Computative Instruction Optimization Pipeline<ul style="list-style-type: none"><li>- Reassociate</li></ul></li></ol> |
|--|

- Constant Propagation
- **Combine Instructions**
- Code Sinking
- SCCP + Simplify CFG
- 5. Loop Optimization Pipelines
  - 5-1. LICM&Rotating&Unswitching
    - Reassociate
    - Loop Simplify
    - Loop-Closed SSA Form(Terminal phi node)
    - Loop Invariant Code Motion
    - Loop Rotate
    - Loop Unswitching
  - 5-2. Loop vectorization
    - Reassociate
    - Loop Simplify
    - Canonicalize induction variables
    - Loop Unroll and Jam + *Vectorize Load & Store*
- 6. Backend Optimization
  - **Liveness Analysis**
  - **Direct Access to Parent Local Registers(Parameter Reduction)**
  - **Register Allocation**
  - LLVM IR to Assembly Static Translation

## II. Basic Information about Custom Passes

This section contains the basic description for custom passes we plan to implement in our project,

### 1. GV to malloc in main()

Naive assembler of our project translates GV into the *malloc* instruction in *main()*. Why not before? Our project specification ensures that there exists *main()* in the IR, allowing this optimization in any cases. This pass will be activated later than the LLVM bundle passes which find and optimize GV uses, but before the heap allocation optimization passes(3-1 from the index above) for maximum efficiency.

### 2. Vectorized Load&Store

\*Highly inspired from the presentation done in 04/28 by Anonymous.

This pass vectorizes the load & store. Our system has 64-bit register, which means that it can load at most 2 *int32*, 4 *int16* registers at once. We may load multiple variables as an *int64* type and then split it using *udiv* & *urem* if a sequential load happens, and vice versa. This reduces the cost of load&store especially when the data type is small. However, sequential memory operations do not happen frequently in normal codes, so we decided to tie this pass with the Loop Unrolling pass(already present in LLVM standard, but need to be modified for compatibility with custom Vectorizing pass).

### 3. Heap to Stack&Argument Lowering

Accessing the heap is not costworthy in most cases, so lowering the variable to stack or register is highly recommended. If the size of the heap allocation is fairly small and the number of any function arguments does not exceed the maximum(16) by this transformation, this pass will bring the heap variable to the stack and send this variable via arguments. Loss caused by increased function arguments will soon be reoptimized by the Interprocedural alloca Sinking and Direct Access to Parent Local Registers passes. Mem2Reg pass should be followed to clean up the *alloca* operations.

c.f. This is the general version of the Local malloc() to alloca pass.

### 4. Interprocedural Alloca Sinking

This pass gets the result of the Heap to Stack&Argument pass, and sinks the *alloca* instruction to the head of the dominating function of all uses, referring to the call graph. This combined with the Dead Argument Elimination reduces the number of function arguments of ancestor functions which is increased by the Heap to Stack&Argument pass.

### 5. Stack/Heap Access grouping

This pass will reorder the memory operations to group stack & heap access as much as it can while not damaging the semantics. In result, we can reduce the number of *reset* operation, and also can save cost for head movement.

### 6. Combine Instructions

This pass will basically be a derivation from *instcombine* pass in LLVM standard. However, project-specific features will be introduced for higher performance. Some potential examples would be:

- `add %x, %x`  $\Rightarrow$  `mul %x, 2`

- and `il %x, %y`  $\Rightarrow$  `mul il %x, %y`
- `add %x, 0`  $\Rightarrow$  `mul %x, 1` (reg-to-reg move instruction)

## 7. Liveness Analysis

Liveness Analysis calculates the live interval of each LLVM IR Registers, and makes into a graph where every LLVM IR registers are a node, and an edge exists between two if their intervals overlap. This graph will be used to determine which register variables cannot be coassigned together via the ‘coloring problem’. Even though unofficial Liveness analysis code exists, for the compatibility of the next pass described we redesign the interface and reimplement it.

## 8. Direct Access to Parent Local Registers

This pass will eventually transform the function parameters to the local register variables. Exploiting the property that `r1~r16` are not discarded but just saved after the function call, functions may freely use the remaining value on the registers. If we use these registers for the arguments, we can save much time in calling & restoring the arguments because to perform operation using arguments we should obviously bring them to the normal registers(`r1~r16`).

The modification is made in the Register Live Interval Graph, which ties two register nodes together. Also, it remembers the change in the function signature so it could be modified later.

c.f. It is pretty obvious that using `argX` register consumes more cost than this method. However, loading arguments by this method may cause register shortage (registers are not enough to compute in a reasonable speed), so ‘rarely accessed read only variables’ can be passed with the `argX` register.

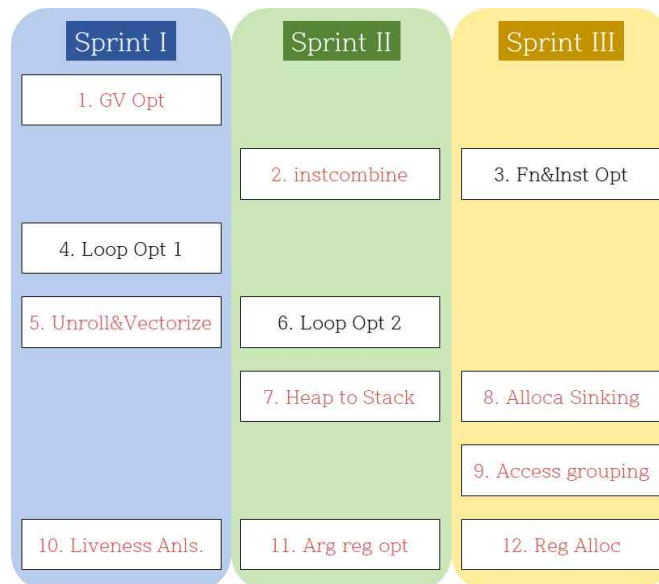
## 9. Register Allocation

These passes are for register allocation. Register Live Interval Graph are colored so that adjacents nodes(IR registers which live intervals overlap) are colored differently. Then they are assigned the physical register in each continuous intervals. The basic procedure is inspired by Chaitin *et al.*

### III. Things to do: 12 Chunks of Tasks and Project Timeline

We divided the whole project into 12 chunks so that workloads are distributed flatly. Task chunks include subtasks as following. Red fonts indicate our custom pass. The dependency between these tasks is described in the following figure. Rows in the diagram indicate the procedure at the left should be completed before the right task for compatibility.

1	Verifying GV optimization + GV to malloc()
2	Modifying Combine Instrurction
3	Verifying Function opt. + Instruction opt.
4	Verifying Loop Opt Pipeline 1
5	Loop unroll + Vectorized load&store
6	Verifying Loop Opt Pipeline 2
7	Heap to Stack&Argument lowering
8	Interprocedural Alloca Sinking
9	Heap/Stack grouping
10	Liveness Analysis + Translation
11	Direct Access to Parent Local Registers
12	Register Allocation



Our team members will each implement and achieve a single task per sprint, total three during the whole project. We distributed these 12 chunks to all members as fairly as possible. Two criteria was applied to this decision.

- When two tasks are dependent to each other, one person should both implement those when possible.
- Estimated workloads should be distributed fairly enough.

Name	Sprint I	Sprint II	Sprint III
임지규	1	2	3
김도현	4	7	8
정덕인	5	6	9
이진우	10	11	12