

# HERCULES: A PATTERN DRIVEN CODE TRANSFORMATION SYSTEM

Christos Kartsaklis

Oscar Hernandez

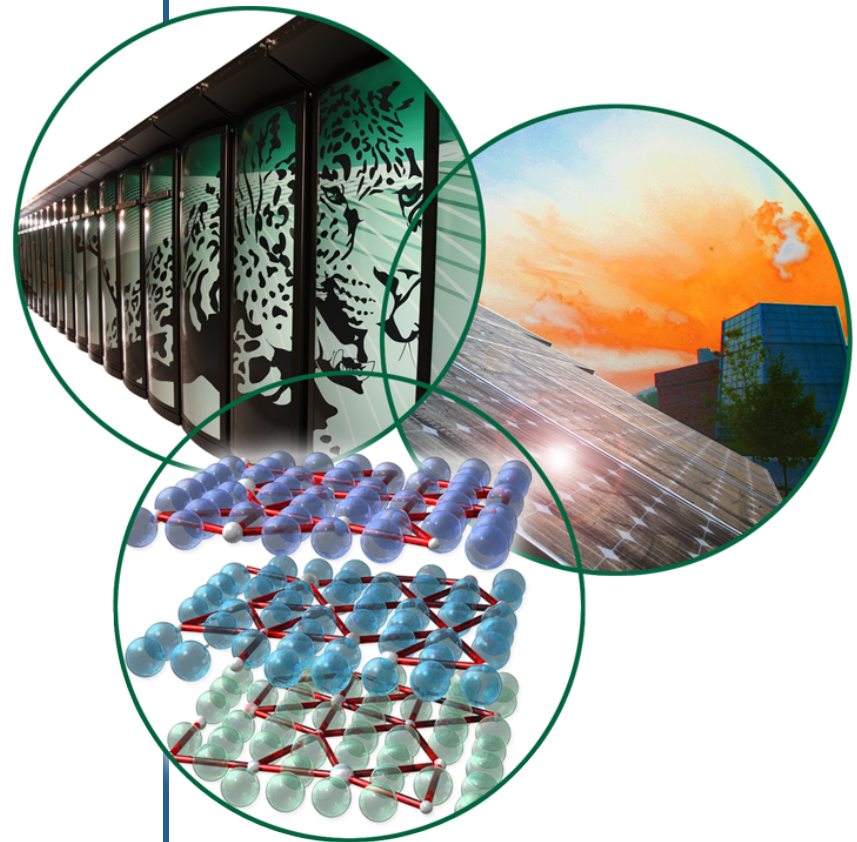
Chung-Hsing Hsu

Thomas Ilsche

Wayne Joubert

Richard Graham

HIPS 2012 (May 21, 2012)

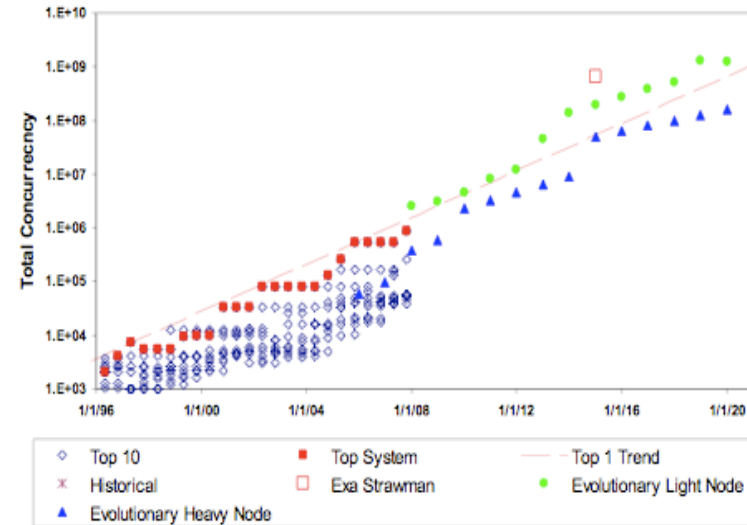


# Take Home Message

**User-oriented tools are more important than ever for programming the new leadership class supercomputers**

# New Leadership Class Machines

- Titan: 10-30PF Cray XK6 (ORNL)
  - Accelerator based
- Path Forward to Exascale:
  - 100,000,000 order cores
  - Extreme levels of parallelism
    - Thread-based programming
    - Task-based execution models
  - Either on die or off die
  - Lower memory footprint per core
  - Deeper memory hierarchies
  - Component failure is the norm



|                        | 2010       | 2018       | Factor Change |
|------------------------|------------|------------|---------------|
| System peak            | 2 Pf/s     | 1 Ef/s     | 500           |
| Power                  | 6 MW       | 20 MW      | 3             |
| System Memory          | 0.3 PB     | 10 PB      | 33            |
| Node Performance       | 0.125 Gf/s | 10 Tf/s    | 80            |
| Node Memory BW         | 25 GB/s    | 400 GB/s   | 16            |
| Node Concurrency       | 12 CPUs    | 1,000 CPUs | 83            |
| Interconnect BW        | 1.5 GB/s   | 50 GB/s    | 33            |
| System Size (nodes)    | 20 K nodes | 1 M nodes  | 50            |
| Total Concurrency      | 225 K      | 1 B        | 4,444         |
| Storage                | 15 PB      | 300 PB     | 20            |
| Input/Output bandwidth | 0.2 TB/s   | 20 TB/s    | 100           |

# Programming Challenges

- Optimization strategies become complex
  - High-levels of concurrency, complex analyses (inter-procedural dependences, scoping of data, global data usage)
- Significant restructuring of applications to new platforms.
  - Significant amount of time consuming, time to invest vs. profitability with **little reuse of knowledge**.
  - Application tied to architecture-specific optimizations
- Multiple programming models and languages.
  - How to map discovered parallelism to programming model, then architecture.
    - What should run on cores, GPUs, across nodes? Load Balance vs. Locality
  - Each with different optimization strategies
- Constant adaptation to new architectures

**Efficient translation of code to architecture**

## Challenges (Cont..)

- Meta programming source-to-source translators rely on backend compilers to generate efficient code:
  - Procedure Inlining
  - Constant Propagation
  - Dead code elimination
  - Common sub-expression elimination
  - Loop optimizations with STL iterators
  - Data flow analysis
  - Alias Analysis

```
template <class T>
struct AxyOp {
    const T * x;
    T * y;
    T alpha, beta;
    void execute(int i)
    { y[i] = alpha*x[i] + beta*y[i]; }
};
```

```
AxyOp<double> op;
op.x = ...; op.alpha = ...;
op.y = ...; op.beta = ...;
node.parallel_for< AxyOp<double> >
    (0, length, op);
```

- Programmers may have to write code in assembler if back-end compiler doesn't generate efficient code.

# Manage the complexity for the user

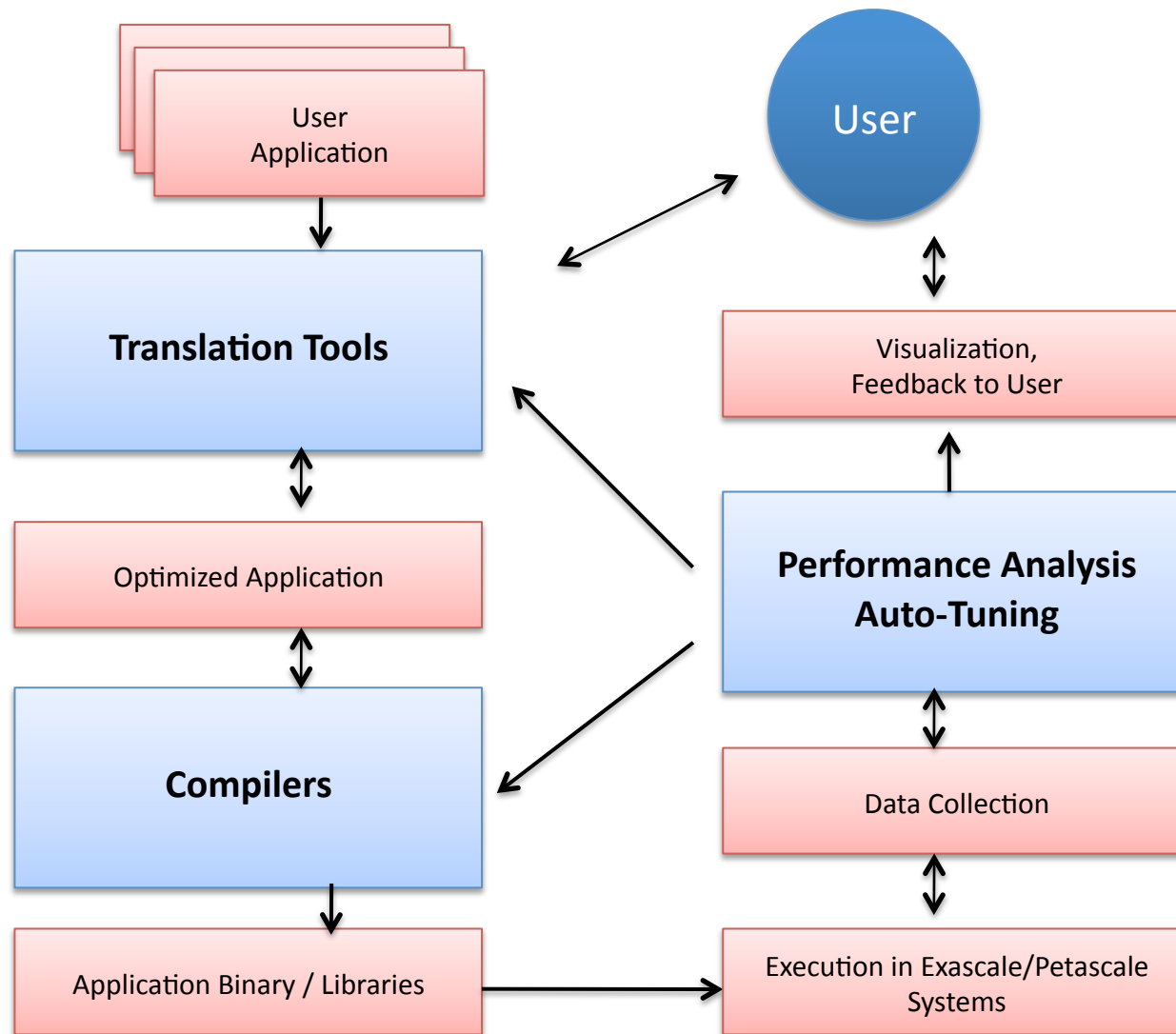
Assist the user in different use-case scenarios:

- Find any nested loop that expands inter-procedurally, where the two outmost loops are parallelizable and inner loops vectorizable with multiply adds.
  - => Parallelize outer loop with OpenMP (auto-scope)
  - => Parallelize second loop with Accelerator Directives (Grid)
  - => Parallelize vector loops with Accelerator Directives (Threadblocks)
- Find vectorizable loops with non-contiguous data accesses and a consumed volume of data less than X
- Find data accesses of a structure inter-procedurally, with no pointers, and that is accessed in loopnests.
  - Re-layout data structure

# Talk Outline

- **Related Work**
- **HERCULES**
- **Lesson Learned**

# Many Stages in Translation





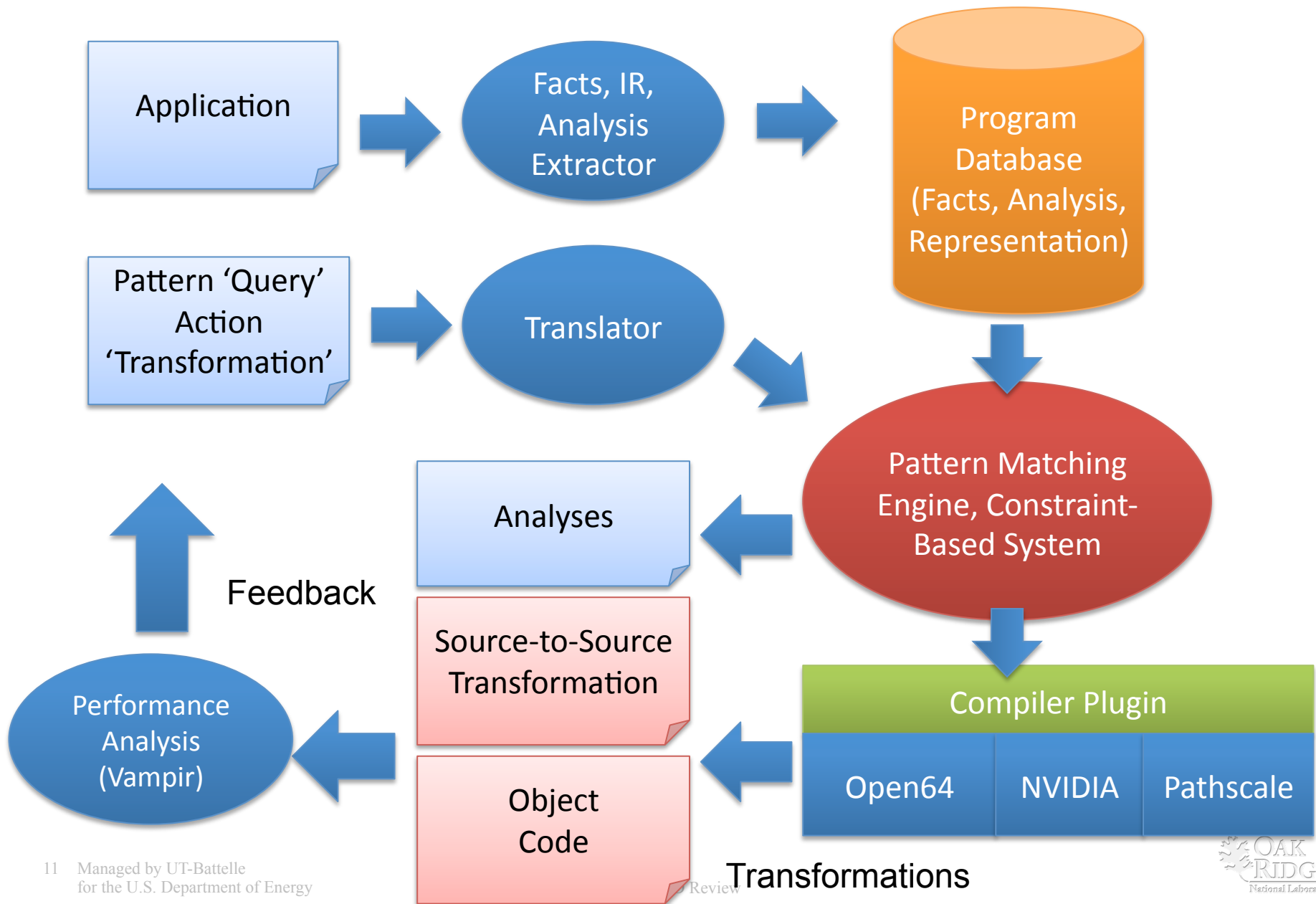
# Related work

- **User-driven source-to-source translators**
  - Poet, Loop Optimizer: Rose (C++/C/Fortran, OpenMP, UPC)
  - ChiLL: Suif (C support)
  - Orio/PluTo: (C/OpenMP)
  - TSF: Forsys (Fortran)
- **Loose integration between translation tool and back-end compiler.**
  - Based on parsers that operate at the abstract syntax tree level, with limited analysis and target-specific transformations.
  - Rely on back-end compiler optimizations
- **Limited feedback from the translation tool to the user.**
- **Strategies tied to a particular user's source code**

# HERCULES

- A pattern driven code translation tool
- Distinctive features:
  - Infrastructure to manage program analysis information for the user to facilitate the understanding of the application
  - Automates the process of applying transformations multiple times throughout the code base
  - Principle of separation of concerns: Application Science vs. Optimizations
  - Documents the transformation process done by computational scientists
  - Reusability of transformation workflow
  - Works with an underlying compiler infrastructure and is a solution that is implementable in compilers.

# Implemented HERCULES Architecture



# Accomplishments

- Working system that can input patterns and transformation scripts and output transformed code and/or binaries.
- Implemented a compiler plug-in infrastructure to enable pattern-matching and transformations at different phases
- Designed a pattern language, and built a directives parser and pattern-matching engine that uses underlying PROLOG technology.
- Progress on output program analysis to the program database
  - Parallelization information, access vectors, cost-models
- Implemented APIs for IR traversal and applying transformations:
  - Loop-transformations, specialization, instrumentation, data transformations, exposed more APIs.
- Worked with applications for requirements of pattern and transformation
  - CAM/SE, Sweep3D, S3D, HPL

# Queries: Pattern Language Definition

- **Incremental approach**

- Define using directives and source code language (C,C++, Fortran)
- Patterns can be generalized incrementally
  - Can easily be used to match a specific source code to be transformed.
- Support for convenience functions
- Can be used to query for analyses
- Pattern can consist of syntactic properties of the code
  - Can be easily extended to support analyses and runtime characteristics of code

Pattern Language Definition:

```
#pragma hercules declare pattern (pattern name, return type)
#pragma hercules symbol bind([variable names],)
#pragma hercules statement bind
#pragma hercules insert ...
#pragma hercules use pattern_name (args)
#pragma hercules bind promote(expre).
```

# Extraction and Pattern Matching Engine

- We translate the program to PROLOG to generate program facts.
  - Mappings between PROLOG facts and compiler intermediate representations are recorded
- Program analysis is also translated to PROLOG
- Pattern translates to PROLOG constraints to be solved with program facts.

Program:

```
#define N 100
int main() {
int a[N], b[N],i
for(i=0; i<N; i++)
    a[i] = b[i] + i;

return a[50]
}
```

IR

```
func_entry
body
  stid i
  intconst 0
do_loop
  le
  ldid I
  intconst 100
body
  stid
  array a[i]
  add
  array b[i]
  ldid i
```

Program  
Facts

(~12x)

```
AST_node(1,0,func_entry,func_entry).
AST_root_of(1,1).
AST_rtype(1,v).
AST_desc(1,v).
AST_kids_of(1,12,[2,3,1,4,5]
AST_has_st(1,1,50,'main').
AST_has_sym(1,12801).
ST_st_idx_to_st(12801,1,50,'main').
ST_index(50,class_func,text).
.....
```

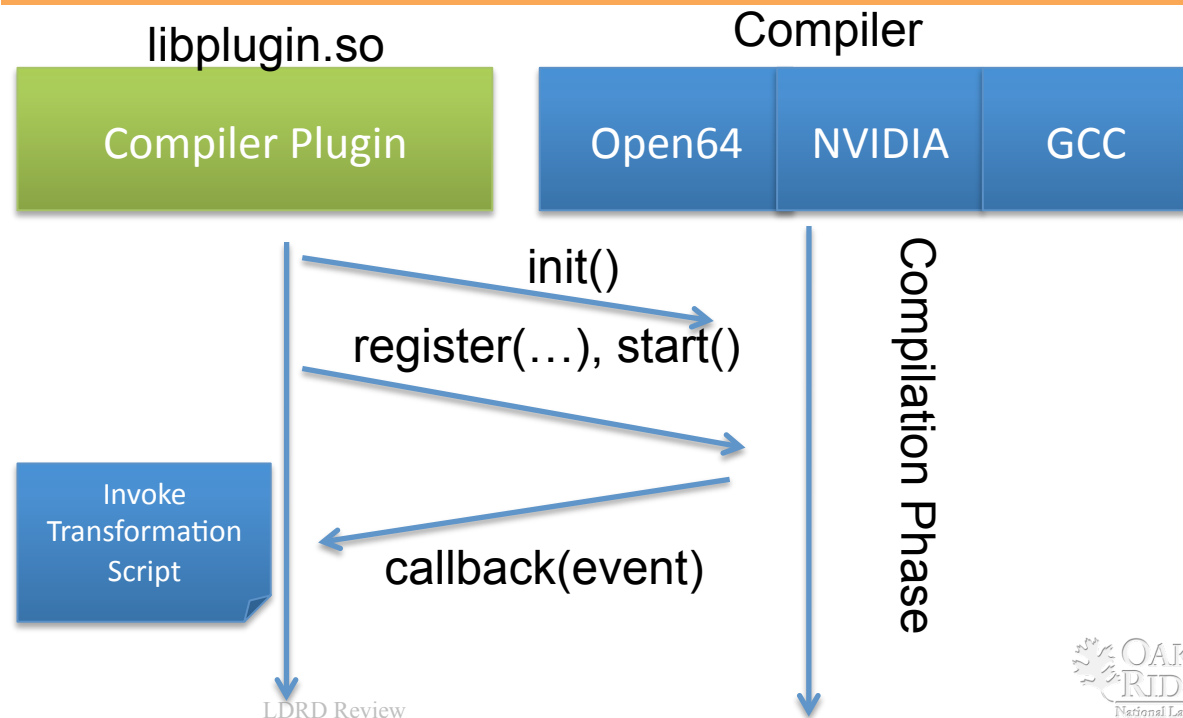
# Compiler Plug-in API

- Implemented a unified compiler plugin-in

```
typedef enum {  
    PLUGIN_REQ_START = 0,  
    PLUGIN_REQ_REGISTER = 1,  
    PLUGIN_REQ_UNREGISTER = 2,  
    PLUGIN_REQ_STATE = 3,  
    PLUGIN_REQ_CURRENT_PHASE = 4,  
    PLUGIN_REQ_STOP = 5,  
    PLUGIN_REQ_PAUSE = 6,  
    PLUGIN_REQ_RESUME = 7,  
    PLUGIN_REQ_LAST = 8  
} PLUGIN_API_REQUEST;
```

```
typedef enum {  
    PLUGIN_EVENT_IPL_BEFORE = 1,  
    PLUGIN_EVENT_IPL_AFTER = 2,  
    PLUGIN_EVENT_IPA_BEFORE = 3,  
    PLUGIN_EVENT_IPA_AFTER = 4,  
    PLUGIN_EVENT_VHO_BEFORE = 5,  
    PLUGIN_EVENT_VHO_AFTER = 6,  
    PLUGIN_EVENT_LNO_BEFORE = 7,  
    PLUGIN_EVENT_LNO_AFTER = 8,  
    PLUGIN_EVENT_WOPT_BEFORE = 9,  
    PLUGIN_EVENT_WOPT_AFTER = 10,  
    PLUGIN_EVENT_CG_BEFORE = 11,  
    PLUGIN_EVENT_CG_AFTER = 12,  
    PLUGIN_EVENT_RESERVED = 13,  
    PLUGIN_EVENT_LAST = 14  
} PLUGIN_API_EVENT;
```

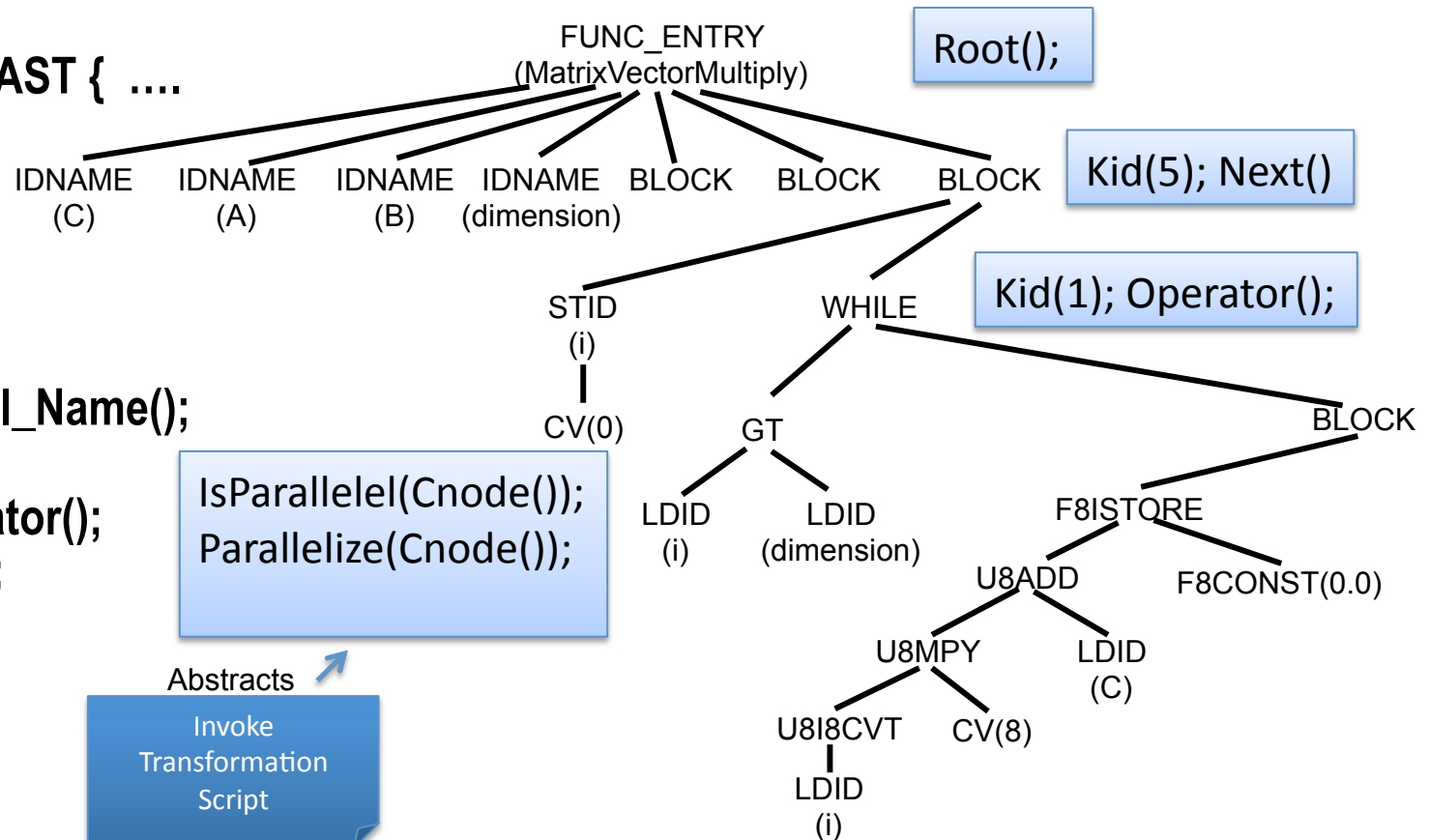
```
class plugin {  
    bool init();  
    void start();  
    void stop();  
    void pause();  
    void resume();  
    void register_event(PLUGIN_API_EVENT e, void (*func)(PLUGIN_API_EVENT e))  
    ....  
}
```



# HERCULES AST & Transformation APIs

```

Class HERCULES_AST { ....
void Root();
void Next();
void Previous();
void Set_Node();
void Kid(INT i);
char * Get_Symbol_Name();
INT Num_Kids();
OPERATOR Operator();
void Statements();
void This_Tree();
void Symbol();
void Type();
void Find();
void Find_Symbols();
void Find_Operator();
void Unroll();
void Specialize();
void Instrument();
void Parent();
void Transformations() };
    
```



IsParallelel(Cnode());  
Parallelize(Cnode());

Abstracts  
Invoke Transformation Script

```

void MatrixVectorMultiply(double *C, double *A, double *B, int dimension) {
int i, j;
for(i=0 ; i<dimension ; i++)
{
C[i] = 0.0;
for(j=0 ; j<dimension ; j++)
C[i] = C[i] + A[i*dimension+j]*B[j];
}
}
    
```



# Example of a Simple Hercules Pattern

```
void mypattern_driver() {  
#pragma hercules pattern declare mpi_loop_pattern (statement FOR, list : statement LEPOINTS)  
#pragma hercules symbol expr1 promote(expression)  
int expr1;  
#pragma hercules statement insert ...  
#pragma hercules symbol i bind(I)  
int i=0;  
#pragma hercules statement bind(FOR) exit_points(LEPOINTS)  
for ( ; i<expr1 ; i++) {  
#pragma hercules statement insert ...  
#pragma hercules pattern use hspl_mpi_callsite(C)  
#pragma hercules statement insert ...  
}  
....  
}
```

Useful for instrumentation transformation or Performance analysis.

Find all loop-nests, each of which contains an MPI call site, return all exit points.

Example of matched code:

```
void foo(int b, int a) {  
int i=0;  
for ( ; i<b ; i++) {  
label1:  
MPI_Send();  
if (a) {  
goto label2;  
} else {  
goto label1;  
}  
if (MPI_Send()) { return; }  
}  
label2:  
bar();  
}
```

# HERCULES in CAM/SE

- Pattern-based parallelization support to reuse transformation logic.

Directive-based pattern definition

## PATTERN DEFINITION:

```
!$hercules pattern declare implicit_camse(statement TARGET)
do ie=nets, nete
!$hercules statement bind TARGET
  do q=1,qsize
    do k=1,nlev
      do j=1,nv
        do l=1,nv
!$hercules statement insert ...
          do i=1,nv
!$hercules statement insert ...
            end do
            divdp4da(l,j,k,q,ie)= rmetdetp(l,j,ie) * ((rdx(ie))*dudx00 ...)
          end do
!$hercules statement insert ...
        end do
      end do
    end do
  end do
end do
end do
end do
```

Transformed

## MATCHED AND TRANSFORMED CODE

```
do ie=nets, nete
!$omp parallel do private(K, Q, J, DUDX00, L, DVDY00I, I),
!$& shared(DVV, METDET, DINV, GRADQ5DA, RMETDETP,
!$& RDX, RDY, IE, DIVDP4DA)
  do q=1,qsize
    do k=1,nlev
      do j=1,nv
        do l=1,nv
          dudx00=0.0d0
          dvd00i=0.0d0
          do i=1,nv
            dudx00 = dudx00 + Dvv(i,l ) * (metdet(i,j,ie)..
            dvd00i = dvd00i + Dvv(i,j ) * (metdet(l,i,ie)) ..
          end do
          divdp4da(l,j,k,q,ie)= rmetdetp(l,j,ie) * ((rdx(ie))*dudx00 ...)
        end do
      end do
    end do
  end do
end do
end do
```

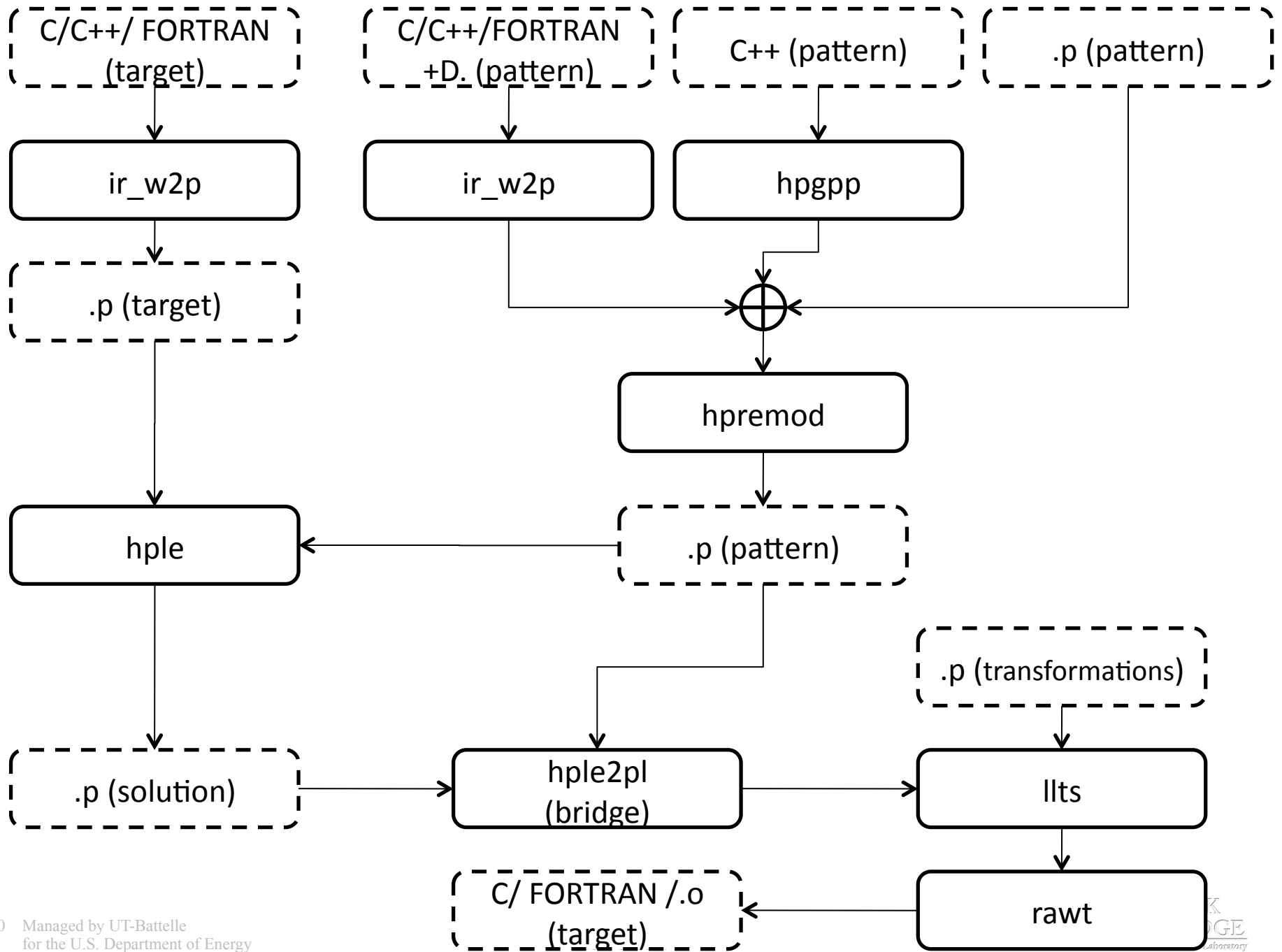
Parallelization support and auto-scoping of variables

## TRANSFORMATION RECIPE:

```
hercules_transformation_for(implicit_camse,1):-
hercules_invoke(s2s_omp_parallelization, arg0)
```

Review





# Lessons Learned

- **Constrained by the choice of the compiler**
  - **Compiler imposes an order/idiom for analysis/transformation.**
  - **Lowering and Normalization.**
- **Influenced by the choice of the pattern language**
  - **Syntax-directed may be limited.**
  - **Generalization versus specialization.**

# THANK YOU

**This research was sponsored by the LDRD program of ORNL, managed by UT-Battelle, LLC for the DOE under Contract DE-AC05-00OR22725**