

Scalable session programming for heterogeneous high-performance systems

Nicholas Ng Nobuko Yoshida Wayne Luk
Imperial College London

September 23, 2013



Outline

Introduction

Session C: Type checking approach

Pabble: MPI Code generation approach

Conclusion and future work



Motivation

- ▶ Parallel architectures
 - ▶ Efficient use of hardware resources
 - ▶ eg. Multicore processors, computer clusters
 - ▶ Difficult to program (correctly)
- ▶ Most common MPI error [Intel survey, SE-HPCS'05]
 - ▶ Communication mismatch (send-receive)
 - ▶ Communication deadlocks



Contributions

Session C Type checking for session primitives

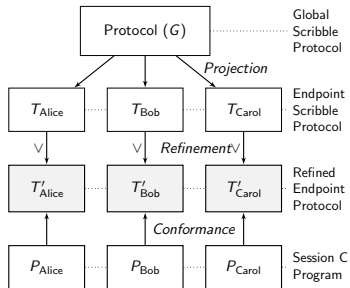
Pabble Session-directed MPI code generation

- ▶ Two approaches to session-based parallel programming
 - ▶ Communication safety
 - ▶ Deadlock freedom
- ▶ Expressing scalable communication topologies as sessions
 - ▶ i.e. in our session-based protocol language



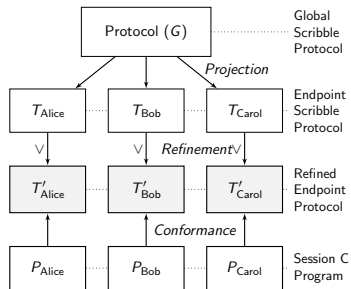
Approach 1 - Session C programming

- ▶ Top-down approach
- ▶ Multiparty session types (MPST)
 - ▶ Communication: duality
 - ▶ Communication safety, deadlock freedom by typing



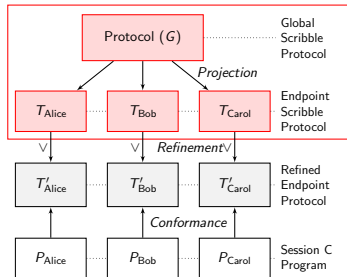
Session C programming: Key reasoning

1. Design protocol in global view
2. Automatic *projection* to endpoint protocol, algorithm preserves safety
3. Write program according to endpoint protocol
4. Check program conforms to protocol
5. \Rightarrow Safe program by design



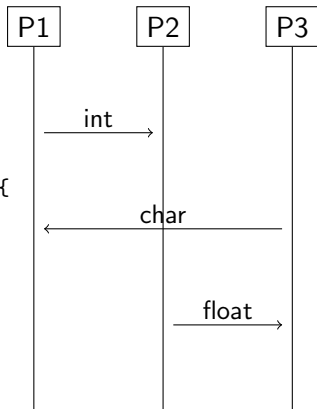
Scribble protocol specification language

- ▶ Communication protocols
- ▶ Interaction by message passing
- ▶ Captures protocol control-flow
- ▶ Developed with industry
 - ▶ Red Hat
 - ▶ Cognizant
 - ▶ Ocean Observatories Initiative
- ▶ More details in tomorrow's demo session



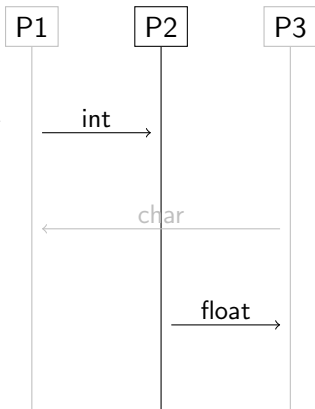
Scribble protocol specification language: Example

```
1  /* Global protocol */
2  protocol Simple
3    (role P1, role P2, role P3) {
4    int  from P1 to P2;
5    char from P3 to P1;
6    float from P2 to P3
7  }
```



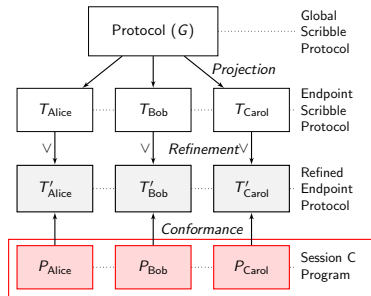
Scribble protocol specification language: Example

```
1  /* Endpoint protocol for P2 */
2  protocol Simple at P2
3    (role P1, role P3) {
4    int from P1;
5
6    float to P3;
7  }
```



Session C runtime

- ▶ Message passing API
- ▶ Built on 0MQ
 - ▶ Fast P2P communication
 - ▶ Lightweight
- ▶ Designed to be simple
 - ▶ Resembles Scribble
 - ▶ Some collective ops support



Session C runtime: Examples

Iteration and message passing

Scribble

```
1  rec X {  
2    (int) to A;  
3    continue X;  
4  }
```

```
1  rec Y {  
2    (int) from B;  
3    continue Y;  
4  }
```

API (simple conditional)

```
1  while (i<3) {  
2    int val = 42;  
3    send_int(&val, 1, A);  
4  }
```

```
1  while (i<3) {  
2    int val; recv_int(&val, 1, B);  
3  }
```



Session C runtime: Examples

Directed choice

Scribble

```
1 choice to B {
2   LABEL0(int) to B;
3 } or {
4   LABEL1(int) to B; }
```

```
1 choice from A {
2   LABEL0(int) from A;
3 } or {
4   LABEL1(int) from A; }
```

API

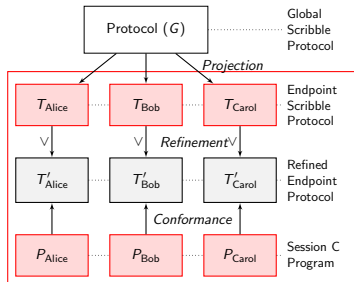
```
1 if (i<3) { // Choice from
2   outbranch(B, LABEL0);
3   send_int(B, 12);
4 } else {
5   outbranch(B, LABEL1);
6   send_char(B, 'A'); }
```

```
1 // Choice to
2 switch (inbranch(A, &label)) {
3   case LABEL0:
4     recv_int(A, &ival); break;
5   case LABEL1:
6     recv_char(A, &cval); break; }
```



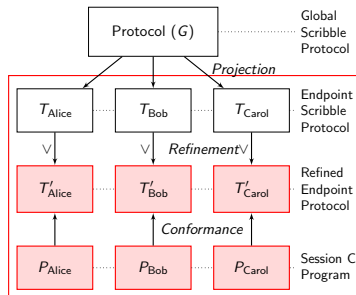
Session Type checking

- ▶ Static analyser
- ▶ Does source code conform to specification?
- ▶ Extract session type from code
 - ▶ Based on usage of API
 - ▶ Based on program flow control
- ▶ Compare w/ endpoint protocol



Session Type checking: Asynchronous optimisation

- ▶ Protocols designed safe
- ▶ Naive impl. inefficient
- ▶ Asynchronous impl.
 - ▶ Non-blocking send
 - ▶ Blocking receive
- ▶ Overlap send/rcv operations
- ▶ Safety by async. subtyping [Mostrous et al., ESOP'09]



Summary (1/2): Session C programming framework

- ▶ Approach: Safety by type checking
- ▶ Protocol-based parallel programming framework
- ▶ Developer friendly Session Types as protocols
- ▶ Implementation with custom API
- ▶ Guarantees communication safety, deadlock free by design



Approach 2: MPI Pabble Code generation approach

- ▶ Scaling: More practical parallel programming
- ▶ Message Passing Interface (MPI) is standard API
- ▶ Associate **Parameterised** MPST with MPI
 - ▶ Type representation (protocol)
 - ▶ Pabble: Parameterised Scribble
 - ▶ Scribble roles with indices
 - ▶ Type check/extraction from source code
 - ▶ Parameterised (dependent) type checking non-trivial
 - ▶ MPI deductive verificationRelated: next talk this session
- ▶ Our solution: Code generation from Pabble protocols



Writing a parallel pipeline in Scribble

```
1 global protocol Ring(role Worker1, role Worker2,  
2   role Worker3, role Worker4) {  
3   rec LOOP {  
4     Data(int) from Worker1 to Worker2 ;  
5     Data(int) from Worker2 to Worker3 ;  
6     Data(int) from Worker3 to Worker4 ;  
7     Data(int) from Worker4 to Worker1 ;  
8     continue LOOP;  
9   }  
10 }
```



Pabble: Parameterised Scribble

- ▶ **Parameterised Scribble** extension
- ▶ Role parameterisation by indices
- ▶ Grouping: Single endpoint protocol for parameterised roles
- ▶ Parametric extension of Scribble
 - ▶ `foreach`, recursion with loop index binding
 - ▶ `if`, conditional execution (multiple roles in single endpoint)
 - ▶ Role index calculation, design based on [Concurrency: state models and Java programs, Magee and Kramer, 2006]
- ▶ Scalable: Supports unbounded number of roles (for some cases)



Indexed interaction statement

Global protocol

```
1 Data(int) from Worker[i:1..9] to Worker[i+1];
```

Endpoint protocol

- ▶ All Workers share an endpoint protocol
- ▶ statements are executed conditionally (by index)

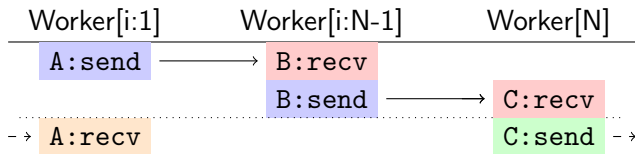
```
1 if Worker[i:2..10] Data(int) from Worker[i-1];
```

```
2 if Worker[i:1..9] Data(int) to Worker[i+1] ;
```



Example: Ring topology in Pabble

```
1 global protocol Ring(role Worker[1..N]) {  
2   rec LOOP {  
3     Data(int) from Worker[i:1..N-1] to Worker[i+1] ;  
4     Data(int) from Worker[N] to Worker[1] ;  
5     continue LOOP;  
6   }  
7 }
```



Ring protocol: Worker endpoint

```
1 local protocol Ring at Worker[1..N](role Worker[1..N]) {
2   rec LOOP {
3     if Worker[i:2..N] Data(int) from Worker[i-1];
4     if Worker[i:1..N-1] Data(int) to Worker[i+1];
5     if Worker[1] Data(int) from Worker[N];
6     if Worker[N] Data(int) to Worker[1];
7     continue LOOP;
8   }
9 }
```



MPI code generation

- ▶ Sessions and MPI: Similar program structure
 - ▶ Pabble also single-source multiple-endpoints
 - ▶ Parameterised role index = MPI ranks
- ▶ Pabble vs. core MPI primitives, e.g.
 - ▶ P2P: Send, Receive
 - ▶ Collective ops: Scatter, Gather, All to All



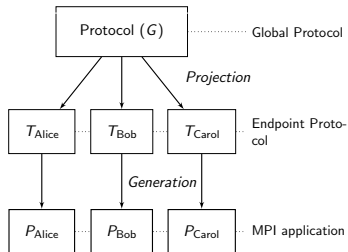
Ring protocol: Simplified MPI code

```
1 MPI_Init(&argc, &argv);
2 MPI_Comm_rank(MPI_COMM_WORLD, &rank);
3 MPI_Comm_size(MPI_COMM_WORLD, &N);
4 while (1) { // rec LOOP
5     // if Worker[i:2..N] Data(int) from Worker[i-1];
6     if (2<=rank && rank<=N) MPI_Recv(.., MPI_INT, rank-1, Data, .. );
7     // if Worker[i:1..N-1] Data(int) to Worker[i+1];
8     if (1<=rank && rank<=N-1) MPI_Send(.., MPI_INT, rank+1, Data, .. );
9     // if Worker[1] Data(int) from Worker[N];
10    if (rank==1) MPI_Recv(.., MPI_INT, N, Data, .. );
11    // if Worker[N] Data(int) to Worker[1];
12    if (rank==N) MPI_Recv(.., MPI_INT, 1, Data, .. );
13 }
14 MPI_Finalize();
```

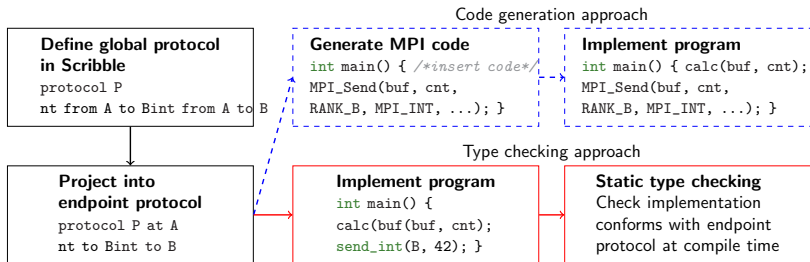


Summary (2/2): MPI code generation from Pabble

- ▶ Approach: Safety by code generation
- ▶ Generate MPI backbone
 - ▶ Communication-correct
- ▶ Pabble indexed roles to rank
- ▶ Supports MPI collective ops



Conclusion: Session-based safe parallel programming



- ▶ Communication safety
- ▶ Deadlock free

Ongoing and future work

- ▶ Extract/verify Session Types from MPI
 - ▶ Huge engineering challenge
 - ▶ Can we infer global types from the endpoint MPI programs?
 - ▶ Ongoing collaboration with FCUL [EuroMPI'12, PLACES'13, BEAT2'13]
- ▶ Applying methodology in different environments
 - ▶ Software-Hardware communication (eg. FPGA, Maxeler)
 - ▶ Heterogeneous hardware (FPGA) code generation via AOP
- ▶ ...



