# Designing Asynchronous Multiparty Protocols with Crash-Stop Failures

Adam D. Barwell[1]    Ping Hou[2]    Nobuko Yoshida[2]    Fangyi Zhou[2,3]

[1]University of St Andrews    [2]University of Oxford    [3]Imperial College London

Stardust Meeting, Glasgow
18 December 2023

University of St Andrews    UNIVERSITY OF OXFORD    Imperial College London

# Concurrent and Communicating Systems

» Are **ubiquitous**
  – Mobile phones, desktops, distributed systems, the internet, *&c.*

» Can be difficult to get **right** unaided
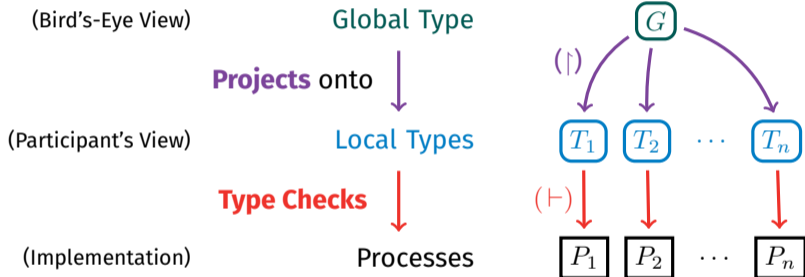  – Communications mismatch, deadlocks, livelocks, *&c.*

# Concurrent and Communicating Systems

» Are ubiquitous
  – Mobile phones, desktops, distributed systems, the internet, *&c.*

» Can be difficult to get right **unaided**
  – Communications mismatch, deadlocks, livelocks, *&c.*

» Some assistance: **Session Types**

# (Multiparty) Session Types

» Formally describe **communications behaviour** between two or more systems

» Communications behaviour is **enforced statically**

# Top-Down Multiparty Session Types

| | | |
|---|---|---|
| (Bird's-Eye View) | Global Type | $G$ |
| | **Projects** onto $(\upharpoonright)$ | |
| (Participant's View) | Local Types | $T_1$ $T_2$ $\cdots$ $T_n$ |
| | **Type Checks** $(\vdash)$ | |
| (Implementation) | Processes | $P_1$ $P_2$ $\cdots$ $P_n$ |

# Top-Down Multiparty Session Types – Advantages

» **Correct-by-construction** behavioural properties
  – Communication safety
  – Deadlock-freedom
  – Liveness

» **Generate** protocol-conforming code
  – Library support in Scala, Rust, Haskell, Erlang, OCaml, &c.

# One Caveat: Fault-Tolerance, or Lack Thereof

» Traditional MPST systems assume a **perfect world**
  – No process failures
  – No message loss, duplication, or corruption

# One Caveat: Fault-Tolerance, or Lack Thereof

» Traditional MPST systems assume a **perfect world**
  - No process failures
  - No message loss, duplication, or corruption

» Distributed systems can, and will, fail
  - Reality is not so convenient…

# One Caveat: Fault-Tolerance, or Lack Thereof

» Traditional MPST systems assume a **perfect world**
  – No process failures
  – No message loss, duplication, or corruption

» Distributed systems can, and will, fail
  – Reality is not so convenient…

» Such MPST systems cannot reason about failures
  – How do participants handle crashes?
  – Do behavioural guarantees still hold in the presence of crashes?
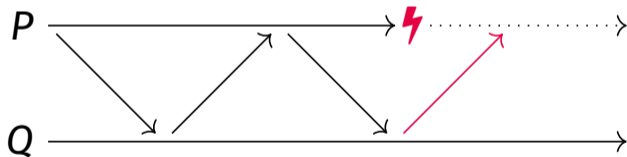
# Asynchronous Multiparty Protocols with Crash-Stop Failures

We present an **asynchronous top-down** MPST theory with **crash-stop failures**.

1. Behavioural Guarantees

2. Optional Reliability Assumptions
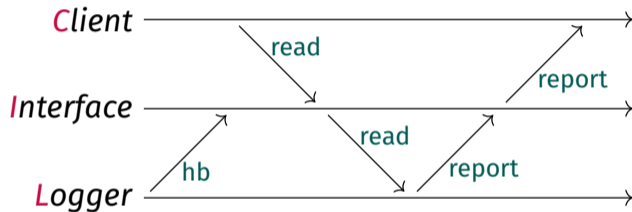
3. Code generation toolchain, TEATRINO

# Crash-Stop Failures

» Processes can **crash arbitrarily**

» Crashed processes **make no progress** and **do not recover**

» Communications channels deliver messages in order and without losses
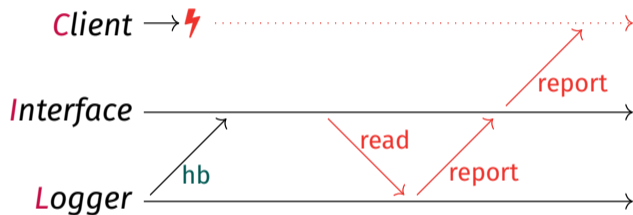
# (Part of) A Simple Distributed Logger

A Client requests the accumulated logs from a distributed Logger via an Interface process.



$G = $ L→I:hb.C→I:read.I→L:read.L→I:report(log).I→C:report(log).end

---

The full version of the protocol (and other variants) can be found in the paper/artefact.
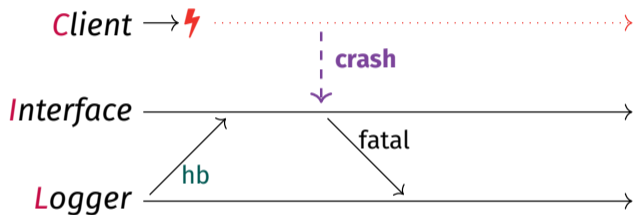
# (Part of) A Simple Distributed Logger

A Client requests the accumulated logs from a distributed Logger via an Interface process.



$$G = \text{L} \rightarrow \text{I:hb}. \ ??$$

# (Part of) A Simple Distributed Logger

A Client requests the accumulated logs from a distributed Logger via an Interface process.



$$G = \text{L→I:hb.C→I:} \left\{ \begin{array}{l} \texttt{read .I→L:read.L→I:report(log).I→C:report(log).end} \\ \texttt{crash .I→L:fatal.end} \end{array} \right\}$$

↑ Reserved label; represents **crash detection**

# Global Type Syntax

$$B \quad ::= \quad \text{int} \mid \text{bool} \mid \text{real} \mid \text{unit} \mid \ldots \qquad \text{Basic types}$$

$$G \quad ::= \quad \text{p} \rightarrow \text{q}^\dagger \colon \{\text{m}_\text{i}(B_i).G_i\}_{i \in I} \qquad\qquad \text{Transmission}$$

$$\mid \quad \text{p}^\dagger \rightsquigarrow \text{q} \colon j \, \{\text{m}_\text{i}(B_i).G_i\}_{i \in I} \; (j \in I) \quad \text{Transmission en route}$$

$$\mid \quad \mu\mathbf{t}.G \qquad\qquad\qquad\qquad\qquad \text{Recursion}$$

$$\mid \quad \mathbf{t} \qquad\qquad\qquad\qquad\qquad\quad\; \text{Type variable}$$

$$\mid \quad \text{end} \qquad\qquad\qquad\qquad\qquad\; \text{Termination}$$

$$\dagger \quad ::= \quad \cdot \quad \mid \quad \lightning \qquad\qquad\qquad\qquad \text{Crash annotation}$$

» Runtime types are not available to the user, only via reduction.

» crash is a reserved label indicating crash-handling branches.

# Crashing Clients

Participant C Crashes.

$$G = \texttt{L}{\rightarrow}\texttt{I:hb.C}{\rightarrow}\texttt{I:} \begin{cases} \texttt{read.I}{\rightarrow}\texttt{L:read.L}{\rightarrow}\texttt{I:report(log).I}{\rightarrow}\texttt{C:report(log).end} \\ \texttt{crash.I}{\rightarrow}\texttt{L:fatal.end} \end{cases}$$

# Crashing Clients

Participant C Crashes.

$$G = \text{L}\rightarrow\text{I:hb.C}\rightarrow\text{I:} \begin{cases} \text{read.I}\rightarrow\text{L:read.L}\rightarrow\text{I:report(log).I}\rightarrow\text{C:report(log).end} \\ \text{crash.I}\rightarrow\text{L:fatal.end} \end{cases}$$

$$\downarrow_{\mathcal{R}}$$

$$G_1 = \text{L}\rightarrow\text{I:hb. } \text{C}^{\frac{1}{2}} \rightsquigarrow \text{I} : \begin{cases} \text{read.I}\rightarrow\text{L:read.L}\rightarrow\text{I:report(log).I}\rightarrow \text{C}^{\frac{1}{2}} :\text{report(log).end} \\ \text{crash.I}\rightarrow\text{L:fatal.end} \end{cases}$$

# Crashing Clients

Participant `C` Crashes.

$$G = \texttt{L} \rightarrow \texttt{I:hb.C} \rightarrow \texttt{I:} \begin{cases} \texttt{read.I} \rightarrow \texttt{L:read.L} \rightarrow \texttt{I:report(log).I} \rightarrow \texttt{C:report(log).end} \\ \texttt{crash.I} \rightarrow \texttt{L:fatal.end} \end{cases}$$

p is unreliable     p occurs in $G$     $G$ is anything but a $\mu$-term

$$\downarrow \quad \quad \frac{ \boxed{p \notin \mathcal{R}} \quad \boxed{p \in \text{roles}(G)} \quad \boxed{G \neq \mu \mathbf{t}.G'} }{ \langle\, \boxed{\mathcal{C}}\, ;G \rangle \xrightarrow{p\,\sharp} _{\mathcal{R}} \langle\, \mathcal{C} \cup \{p\};\ \boxed{G \sharp p}\, \rangle } \; [\text{GR-}\sharp]$$

Set of crashed roles      'Remove' p from $G$

$$G_1 = \texttt{L} \rightarrow \texttt{I:hb.}\; \boxed{\texttt{C}^{\sharp} \rightsquigarrow \texttt{I}} : \begin{cases} \texttt{read.I} \rightarrow \texttt{L:read.L} \rightarrow \texttt{I:report(log).I} \rightarrow \boxed{\texttt{C}^{\sharp}} \texttt{:report(log).end} \\ \texttt{crash.I} \rightarrow \texttt{L:fatal.end} \end{cases}$$

# But Some are More Reliable than Others

| **No** processes can crash | ← But what about **some** processes can crash? → | **All** processes can crash |
| --- | --- | --- |

» Some participants may represent services that can be **assumed reliable**

# But Some are More Reliable than Others



» Some participants may represent services that can be **assumed reliable**

## Optional Reliability Assumptions

- » Specific participants can be marked as **reliable**
- » Reliable participants **will not crash**
- » Reduction rules, *&c.* conform to these assumptions

# Crashing Clients

Participant I Detects

$$G_1 = \text{L}\rightarrow\text{I:hb.C}^{\xi} \rightsquigarrow \text{I:} \begin{cases} \text{read.I}\rightarrow\text{L:read.L}\rightarrow\text{I:report(log).I}\rightarrow\text{C}^{\xi}\text{:report(log).end} \\ \text{crash.I}\rightarrow\text{L:fatal.end} \end{cases}$$

$$\downarrow_{\mathcal{R}}^{*}$$

# Crashing Clients

Participant I Detects

$$G_1 = \texttt{L}{\rightarrow}\texttt{I}\texttt{:hb.C}^{\lightning} \rightsquigarrow \texttt{I:} \begin{cases} \texttt{read.I}{\rightarrow}\texttt{L:read.L}{\rightarrow}\texttt{I:report(log).I}{\rightarrow}\texttt{C}^{\lightning}\texttt{:report(log).end} \\ \texttt{crash.I}{\rightarrow}\texttt{L:fatal.end} \end{cases}$$

$$\downarrow_{\mathcal{R}}^{*}$$

$$G_4 = \texttt{C}^{\lightning} \rightsquigarrow \texttt{I:} \begin{cases} \texttt{read.I}{\rightarrow}\texttt{L:read.L}{\rightarrow}\texttt{I:report(log).I}{\rightarrow}\texttt{C}^{\lightning}\texttt{:report(log).end} \\ \texttt{crash.I}{\rightarrow}\texttt{L:fatal.end} \end{cases}$$

# Crashing Clients

Participant I Detects

$$G_4 = C^{\maltese} \rightsquigarrow I: \begin{cases} \texttt{read.I} \rightarrow \texttt{L:read.L} \rightarrow \texttt{I:report(log).I} \rightarrow C^{\maltese}: \texttt{report(log).end} \\ \texttt{crash.I} \rightarrow \texttt{L:fatal.end} \end{cases}$$

$$\downarrow_{\mathcal{R}}$$

$$G_5 = I \rightarrow \texttt{L:fatal.end}$$

# Crashing Clients

$$G_4 = \mathtt{C}^{\sharp} \rightsquigarrow \mathtt{I} : \begin{cases} \mathtt{read.I} \to \mathtt{L:read.L} \to \mathtt{I:report(log).I} \to \mathtt{C}^{\sharp} : \mathtt{report(log).end} \\ \mathtt{crash.I} \to \mathtt{L:fatal.end} \end{cases}$$

$$\downarrow_{\mathcal{R}} \quad \frac{j \in I \quad \mathtt{m}_j = \mathsf{crash}}{\langle \mathcal{C}; \mathtt{p}^{\sharp} \rightsquigarrow \mathtt{q}{:} j \left\{ \mathtt{m}_\mathtt{i}(B_i).G_i' \right\}_{i \in I} \rangle \xrightarrow{\mathtt{q} \odot \mathtt{p}}_{\mathcal{R}} \langle \mathcal{C}; G_j' \rangle} \; \text{[GR-}\odot\text{]}$$

$$G_5 = \mathtt{I} \to \mathtt{L:fatal.end}$$

# Local Type Syntax

$$B \quad ::= \quad \text{int} \mid \text{bool} \mid \text{real} \mid \text{unit} \mid \ldots \quad \text{Basic types}$$

$$
\begin{aligned}
S, T \quad ::= \quad & \text{p}\&\{\text{m}_\text{i}(B_i).T_i\}_{i\in I} && \text{External choice} \\
\mid \quad & \text{p}\oplus\{\text{m}_\text{i}(B_i).T_i\}_{i\in I} && \text{Internal choice} \\
\mid \quad & \mu\textbf{t}.T && \text{Recursion} \\
\mid \quad & \textbf{t} && \text{Type variable} \\
\mid \quad & \text{end} && \text{Termination} \\
\mid \quad & \boxed{\text{stop}} && \text{Crash}
\end{aligned}
$$

» The  runtime  stop type is the type of crashed participants

» crash is a reserved label indicating crash-handling branches

# Type Checking Distributed Logging with Local Types

$$T_\text{C} \;=\; \text{I} \oplus \text{read.I} \& \text{report(log).end}$$

$$P_\text{C} \;=\; \text{I!read.I?report}(x).\mathbf{0}$$

$$T_\text{I} \;=\; \text{L\&hb.C\&} \begin{Bmatrix} \text{read.L} \oplus \text{read.L\&report(log).C} \oplus \text{report(log).end} \\ \text{crash.L} \oplus \text{fatal.end} \end{Bmatrix}$$

$$P_\text{I} \;=\; \text{L?hb.} \sum \begin{Bmatrix} \text{C?read.L!read.L?report}(x).\text{C!report}\langle x\rangle.\mathbf{0} \\ \text{C?crash.L!fatal.}\mathbf{0} \end{Bmatrix}$$

---

We largely elide the session $\pi$-calculus here, but full details can be found in the paper.

# Local Type Context Reduction

$$T_{\mathrm{C}} \;=\; \mathtt{I}\oplus\mathrm{read.I\&report(log).end} \qquad\qquad T_{\mathrm{C}} \;=\; \mathrm{stop}$$
$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \rightarrow$$
$$P_{\mathrm{C}} \;=\; \mathtt{I!read.I?report}(x).\mathbf{0} \qquad\qquad\qquad P_{\mathrm{C}} \;=\; \lightning$$

» Rules are largely standard
  – Crashing, crash-detection, &c. rules are novel

» Queues handle asynchronous message passing
  – Queues are made unavailable when the corresponding participant has crashed

# Projecting Global Types into Local Types

Relates *G* and *T* for a given $p$.

$$(q \to r^\dagger : \{m_i(B_i).G_i\}_{i \in I}) \upharpoonright_{\mathcal{R}} p = \begin{cases} r \oplus \{m_i(B_i).(G_i \upharpoonright_{\mathcal{R}} p)\}_{i \in \{j \in I \,|\, m_j \neq \text{crash}\}} & \text{if } p = q \\ q \& \{m_i(B_i).(G_i \upharpoonright_{\mathcal{R}} p)\}_{i \in I} & \text{if } p = r, \text{ and } q \notin \mathcal{R} \text{ implies} \\ & \exists k \in I : m_k = \text{crash} \\ \prod_{i \in I} G_i \upharpoonright_{\mathcal{R}} p & \text{if } p \neq q, \text{ and } p \neq r \end{cases}$$

$$(q^\dagger \rightsquigarrow r : j \{m_i(B_i).G_i\}_{i \in I}) \upharpoonright_{\mathcal{R}} p = \begin{cases} G_j \upharpoonright_{\mathcal{R}} p & \text{if } p = q \\ q \& \{m_i(B_i).(G_i \upharpoonright_{\mathcal{R}} p)\}_{i \in I} & \text{if } p = r, \text{ and } q \notin \mathcal{R} \text{ implies} \\ & \exists k \in I : m_k = \text{crash} \\ \prod_{i \in I} G_i \upharpoonright_{\mathcal{R}} p & \text{if } p \neq q, \text{ and } p \neq r \end{cases}$$

$$(\mu \mathbf{t}.G) \upharpoonright_{\mathcal{R}} p = \begin{cases} \mu \mathbf{t}.(G \upharpoonright_{\mathcal{R}} p) & \text{if } p \in G \text{ or } \text{fv}(\mu \mathbf{t}.G) \neq \emptyset \\ \text{end} & \text{otherwise} \end{cases} \qquad \mathbf{t} \upharpoonright_{\mathcal{R}} p = \mathbf{t}$$

$$\text{end} \upharpoonright_{\mathcal{R}} p = \text{end}$$

---

The (full) merging operator definition can be found in the paper.

# Projection Begat Properties

1. Association of Global Type and Configuration (Soundness and Completeness)
   - Global and local types do not diverge when reducing

---

Formal definitions and proofs can be found in the paper.

# Projection Begat Properties

1. Association of Global Type and Configuration (Soundness and Completeness)
   - Global and local types do not diverge when reducing

2. Configuration Safety
   - There are no label mismatches
   - Each receiver must be able to handle the potential crash of the (unreliable) sender

---

Formal definitions and proofs can be found in the paper.

# Projection Begat Properties

1. Association of Global Type and Configuration (Soundness and Completeness)
   - Global and local types do not diverge when reducing

2. Configuration Safety
   - There are no label mismatches
   - Each receiver must be able to handle the potential crash of the (unreliable) sender

3. Deadlock-Freedom (Progress)
   - Local types are able to reduce until they terminate or crash

Formal definitions and proofs can be found in the paper.

# Projection Begat Properties

1. Association of Global Type and Configuration (Soundness and Completeness)
   - Global and local types do not diverge when reducing

2. Configuration Safety
   - There are no label mismatches
   - Each receiver must be able to handle the potential crash of the (unreliable) sender

3. Deadlock-Freedom (Progress)
   - Local types are able to reduce until they terminate or crash

4. Liveness
   - Every pending internal/external choice is eventually triggered (by message transmission or crash detection)

---

Formal definitions and proofs can be found in the paper.

# Asynchronous Multiparty Protocols with Crash-Stop Failures

We present an **asynchronous top-down** MPST theory with **crash-stop failures**.

✓ Behavioural Guarantees

✓ Optional Reliability Assumptions

⇒ Code generation toolchain, TEATRINO

# TEATRINO



» SCRIBBLE-inspired **code generation toolchain**
- Consumes SCRIBBLE protocols
- Produces protocol-conforming SCALA code using the EFFPI concurrency library
- Generated code is executable

» **Implements** our (non-runtime) Global and Local Types, and Projection

» **Extends** both the EFFPI and SCRIBBLE syntaxes with crash-stop failures

# The (Extended) SCRIBBLE Protocol Description Language

» We support (a less sugary subset of) the version accepted by $\nu$**SCR**[1]
  – No support for `do`-notation – recursion is expressed via `rec` and `continue`
  – No support for auxiliary protocol definitions

## Extensions

1. `reliable` role declarations

2. Reserved `crash` label

---

[1]https://nuscr.dev/

# The (Extended) SCRIBBLE Protocol Description Language

$$G = \text{L}\rightarrow\text{I:hb.C}\rightarrow\text{I:} \begin{cases} \text{read.I}\rightarrow\text{L:read.L}\rightarrow\text{I:report(log).I}\rightarrow\text{C:report(log).end} \\ \text{crash.I}\rightarrow\text{L:fatal.end} \end{cases}$$

```
global protocol G(reliable role L, role C, reliable role I) {
  hb from L to I;
  choice at C {
    read from C to I;
    read from I to L;
    report(Log) from L to I;
    report(Log) from I to C;
  } or {
    crash from C to I;
    fatal from I to L;
  }
}
```

# The (Extended) EFFPI Concurrency Library

» Embedded Domain Specific Language for SCALA 3
» Leverages type features to **represent local types directly in code**
  – Union types, match types, and dependent and polymorphic function types

## Extensions

1. Support for crash-handling branches
   – New type-level receive construct: `InErr`
   – New value-level receive construct: `receiveErr`

2. Support for crash detection
   – Implemented using timeouts

---

Original version of EFFPI: `https://github.com/alcestes/effpi`

# The (Extended) EFFPI Concurrency Library

$$T_{\text{I}} = \text{L\&hb.C\&} \begin{cases} \text{read.L} \oplus \text{read.L\&report(log).C} \oplus \text{report(log).end} \\ \text{crash.L} \oplus \text{fatal.end} \end{cases}$$

```
type I[C0 <: InChan[Hb],
       C1 <: OutChan[Fatal],
       C2 <: InChan[Read],
       C3 <: InChan[Report],
       C4 <: OutChan[Report]] =
  In[C0, Hb, (X <: Hb) =>
    InErr[C2, Read,
      (Y <: Read) =>
        Out[C3,Read] >>: In[C4, Report, (Z <: Log) => Out[C5, Report]],
      (Err <: Throwable) => Out[C2,Fatal]
  ]]
```
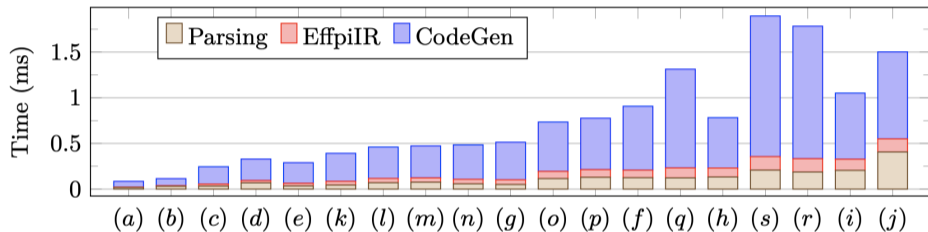
# Evaluating TEATRINO

## Expressiveness

- » Applied to examples from session type and distributed systems literature
- » Standard examples extended with crash-handling behaviour
- » Two patterns: **graceful failure** and **failover**

## Feasibility

- » We give generation times for all of our examples
- » We report times for the three main generation phases:
  1. Parsing
  2. EffpiIR Generation
  3. Code Generation

# Evaluating TEATRINO



» Evaluated on 19 protocols taken from session type and distributed system literature
  – Code generation times all under 3 milliseconds

# Summary

» Asynchronous top-down MPST theory with crash-stop failures:
  – Support for fully-reliable to fully-unreliable protocols
  – Safety, deadlock-freedom, and liveness guarantees

» TEATRINO: toolchain support for generating protocol-conformant SCALA code

» Future work:
  – Investigate different crash and failure models (e.g. crash-recover, link failures)

## Links

» Full version: `https://arxiv.org/abs/2305.06238`

» Artefact: `https://doi.org/10.5281/zenodo.7714132`

» Artefact Source: `https://github.com/adbarwell/ECOOP23-Artefact`