# A Theory for Observational Fault Tolerance (Extended Abstract)

Adrian Francalanza[1] and Matthew Hennessy[2]

[1] University of Malta, Msida MSD 06, Malta, `afra1@um.edu.mt`
[2] University of Sussex, Brighton BN1 9RH, England, `matthewh@sussex.ac.uk`

**Abstract.** In general, faults cannot be prevented; instead, they need to be tolerated to guarantee certain degrees of software dependability. We develop a theory for fault tolerance for a distributed pi-calculus, whereby locations act as units of failure and redundancy is distributed across independently failing locations. We give formal definitions for fault tolerant programs in our calculus, based on the well studied notion of contextual equivalence. We then develop bisimulation proof techniques to verify fault tolerance properties of distributed programs and show they are sound with respect to our definitions for fault tolerance.

## 1 Introduction

One reason for the study of programs in the presence of *faults*, i.e. defects at the lowest level of abstractions [2], is to be able to construct more *dependable* systems, meaning systems exhibiting a high probability of *behaving* according to their *specification* [13]. System dependability is often expressed through attributes like maintainability, availability, safety and *reliability*, the latter of which is defined as a measure of the *continuous delivery* of *correct behaviour*, [13]. There are a number of approaches for achieving system dependability in the presence of faults, ranging from fault removal, fault prevention and *fault tolerance*.

The fault tolerant approach to system dependability consist of various techniques that employ *redundancy* to prevent faults from generating *failure*, i.e. abnormal behaviour caused by faults [2]. Two forms of redundancy are *space redundancy* (*replication*), i.e. using several copies of the same system components, and *time redundancy*, i.e. performing the same chunk of computation more than once. Certain fault tolerant techniques are based on *fault detection* which subsequently trigger *fault recovery*. If enough redundancy is used, fault recovery can lead to *fault masking*, where the specified behaviour is preserved without noticeable glitch.

Fault tolerance is of particular relevance in distributed computing; distribution yield a natural notion of *partial failure*, whereby faults affect a *subset* of the computation. Partial failure, in turn, gives scope for introducing redundancy as *replication*, distributed across independently failing entities such as locations. In general, the higher the replication, the greater the potential for fault tolerance. Nevertheless, fault tolerance also depends on how replicas are managed. One classification, due to [13], identifies three classes, namely *active replication* (all replicas are invoked for every operation), *passive replication* (operations are invoked on primary replicas and secondary replicas are

updated in batches at checkpoints), and *lazy replication* (a hybrid of the previous two, exploiting the separation between write and read operations).

In this paper we address fault tolerance in a distributed setting, focussing on simple examples using *stateless* (read-only) replicas which are invoked only once. We code these examples in $D\pi$ [8] with failing locations [5], a simple distributed version of the standard $\pi$-calculus [11], where the locations that host processes model closely physical network nodes.

*Example 1.* Consider the systems $\mathsf{server}_i$, three server implementations accepting client requests on channel *req* with two arguments, $x$ being the value to process and $y$ being the reply channel on which the answer is returned.

$$\mathsf{server}_1 \;\Leftarrow\; (\nu\,data) \left( \begin{array}{l} l[\![ req?(x,y).\mathsf{go}\,k_1.data!\langle x,y,l\rangle ]\!] \\ |\; k_1[\![ data?(x,y,z).\mathsf{go}\,z.y!\langle f(x)\rangle ]\!] \end{array} \right)$$

$$\mathsf{server}_2 \;\Leftarrow\; (\nu\,data) \left( \begin{array}{l} l \left[\!\!\left[ req?(x,y).(\nu sync) \left( \begin{array}{l} \mathsf{go}\,k_1.data!\langle x, sync, l\rangle \\ |\; \mathsf{go}\,k_2.data!\langle x, sync, l\rangle \\ |\; sync?(x).y!\langle x\rangle \end{array} \right) \right]\!\!\right] \\ |\; k_1[\![ data?(x,y,z).\mathsf{go}\,z.y!\langle f(x)\rangle ]\!] \\ |\; k_2[\![ data?(x,y,z).\mathsf{go}\,z.y!\langle f(x)\rangle ]\!] \end{array} \right)$$

$$\mathsf{server}_3 \;\Leftarrow\; (\nu\,data) \left( \begin{array}{l} l \left[\!\!\left[ req?(x,y).(\nu sync) \left( \begin{array}{l} \mathsf{go}\,k_1.data!\langle x, sync, l\rangle \\ |\; \mathsf{go}\,k_2.data!\langle x, sync, l\rangle \\ |\; \mathsf{go}\,k_3.data!\langle x, sync, l\rangle \\ |\; sync?(x).y!\langle x\rangle \end{array} \right) \right]\!\!\right] \\ |\; k_1[\![ data?(x,y,z).\mathsf{go}\,z.y!\langle f(x)\rangle ]\!] \\ |\; k_2[\![ data?(x,y,z).\mathsf{go}\,z.y!\langle f(x)\rangle ]\!] \\ |\; k_3[\![ data?(x,y,z).\mathsf{go}\,z.y!\langle f(x)\rangle ]\!] \end{array} \right)$$

Requests are forwarded to internal databases, denoted by the scoped channel *data*, distributed and replicated across the auxiliary locations $k_i$. A database looks up the mapping of the value $x$ using some unspecified function $f(-)$ and returns the answer, $f(x)$, back on port $y$. When multiple replicas are used, as in $\mathsf{server}_{2,3}$, requests are sent to all replicas in an arbitrary fashion, without the use of failure detection, and multiple answers are synchronised at $l$ on the scoped channel *sync*, returning the first answer received on $y$.

The theory developed in [5] enables us to *differentiate* between these systems, based on the different behaviour observed when composed with systems such as

$$\mathsf{client} \;\Leftarrow\; l[\![ req!\langle v, ret\rangle ]\!]$$

in a setting where locations may fail. Here we go one step further, allowing us to *quantify* in some sense the difference between these systems. Intuitively, if locations $k_i, i = 1, 2, 3$, can fail in fail-stop fashion[12] and observations are limited to location $l$ only, then $\mathsf{server}_2$ seems to be more *fault tolerant* than $\mathsf{server}_1$; observers limited to $l$, such as $\mathsf{client}$, cannot observe changes in behaviour in $\mathsf{server}_2$ when *at most 1* location

from $k_i$ fails. Similarly, server$_3$ is more *fault tolerant* than server$_1$ and server$_2$ because server$_3$ | client preserves its behaviour at *l up to 2* faults occurring at $k_{1..3}$.

In this paper we give a formal definition of when a system is deemed to be fault tolerant up to *n*-faults, which coincides with this intuition. As in [5] we need to consider systems $M$, running on some network, which we will represent as $\Gamma \triangleright M$. Then we will say that $M$ is fault-tolerant up to *n* faults if

$$F^n[\Gamma \triangleright M] \cong \Gamma \triangleright M \qquad (1)$$

where $F^n[\ ]$ is some context which induces at most *n* faults, and $\cong$ is some behavioural equivalence between systems descriptions. A key aspect of this behavioural equivalence is the implicit separation between *reliable* locations, which are assumed not to fail, and *unreliable* locations, which may fail. In the above example *l* is reliable, at which observations can be made, while the $k_i$ are assumed unreliable, subject to failure. Furthermore it is essential that observers not have access to to these unreliable locations, at any time during a computation. Otherwise (1) would no longer represent $M$ being fault tolerant; for example we would no longer have

$$F^1[\Gamma \triangleright \text{server}_2] \cong \Gamma \triangleright \textbf{server}_2$$

as an observer with access to $k_i$ would be able to detect possible failures in $F^1[\Gamma \triangleright$ server$_2]$, not present in $\Gamma \triangleright$ **server**$_2$.

We enforce this separation between reliable, observable, locations, and unreliable, unobservable, locations, using a simple type system in which the former are called *public*, and the latter *confined*. This is outlined in Section 2, where we also formally define the language we use, DπLoc, give its reduction semantics, and also outline the behavioural equivalence $\cong$; this last is simply an instance of *reduction barbed congruence*, [6], modified so that observations can only be made at public locations. In Section 3 we give our formal definition of fault-tolerance; actually we give two versions of (1) above, called *static* and *dynamic* fault tolerance; we also motivate the difference with examples. Proof techniques for establishing fault tolerance are given in Section 4; in particular we give a complete co-induction characterisation of $\cong$, using labelled actions, and some useful up-to techniques for presenting witness bisimulations. In Section 5 we refine these proof techniques for the more demanding fault tolerant definition, *dynamic fault tolerance*, using *simulations*. Finally Section 6 outlines the main contributions of the paper and discusses future and related work.

## 2   The Language

We assume a set of *variables* Vars, ranged over by $x, y, z, \ldots$ and a separate set of *names*, Names, ranged over by $n, m, \ldots$, which is divided into locations, Locs, ranged over by $l, k, \ldots$ and channels, Chans, ranged over by $a, b, c, \ldots$. Finally we use $u, v, \ldots$ to range over the set of *identifiers*, consisting of either variables and names.

The syntax of our language, DπLoc, is a variation of Dπ [8] and is given in Table 1. The main syntactic category is that of *systems*, ranged over by $M, N$: these are essentially a collection of *located processes*, or *agents*, composed in parallel where location

Table 1. *Syntax of typed DπF*

**Types**

| | | | | | |
|---|---|---|---|---|---|
| $T ::= ch_v\langle\tilde{P}\rangle \mid loc_v^s$ | *(stateful types)* | | $s ::= a \mid d$ | *(status)* |
| $U ::= ch_v\langle\tilde{P}\rangle \mid loc_v$ | *(stateless types)* | | $v ::= p \mid c$ | *(visibility)* |
| $P ::= ch_p\langle\tilde{P}\rangle \mid loc_p$ | *(public stateless types)* | | | |

**Processes**

| | | | | |
|---|---|---|---|---|
| $P, Q ::= u!\langle V\rangle.P$ | *(output)* | | $\mid u?(X).P$ | *(input)* |
| $\mid$ if $v = u$ then $P$ else $Q$ | *(matching)* | | $\mid * u?(X).P$ | *(replicated input)* |
| $\mid (\nu n : T)P$ | *(channel/location definition)* | | $\mid$ go $u.P$ | *(migration)* |
| $\mid \mathbf{0}$ | *(inertion)* | | $\mid P\mid Q$ | *(fork)* |
| $\mid$ ping $u.P$ else $Q$ | *(status testing)* | | | |

**Systems**

| | | | | |
|---|---|---|---|---|
| $M, N, O ::= l[\![P]\!]$ | *(located process)* | | $\mid N\mid M$ | *(parallel)* |
| $\mid (\nu n : T)N$ | *(hiding)* | | | |

and channel names may be scoped to a subset of agents. The syntax for processes, $P, Q$, is an extension of that in Dπ: there is input and output on channels - here $V$ is a tuple of identifiers, and $X$ a tuple of variables, to be interpreted as a pattern - and standard forms of parallel composition, inertion, replicated input, local declarations, a test for equality between identifiers and migration. The only addition on the original Dπ is ping $k.P$ else $Q$, which tests for the *status* of $k$ in the style of [1, 10] and branches to $P$ if $k$ is alive and $Q$ otherwise. For these terms we assume the standard notions of *free* and *bound* occurrences of both names and variables, together with the associated concepts of $\alpha$-conversion and *substitution*. We also assume that systems are *closed*, that is they have no free variable occurrences.

As explained in the Introduction we use a variation (and simplification) of the type system of Dπ [8] in which the the two main categories, channels and locations, are now annotated by visibility constraints, giving $ch_v\langle\tilde{P}\rangle$ and $loc_v$, where $v$ may either be p (i.e. public) or c (i.e. confined); in Table 1 these are called *stateless* types, and are ranged over by U. As explained in [5] a simple reduction semantics can be defined if we also allow types which record the status of a location, whether it is alive, a, or dead, d; these are refereed to as *stateful* types, ranged over by T. Finally P ranges over *public* types, the types assigned to all names which are visible to observers.

*Type System:* $\Gamma$ denotes a type environment, an unordered list of tuples assigning a single *stateful* type to names, and we write $\Gamma \vdash n : T$ to mean that $\Gamma$ assigns the type T to $n$; when it is not relevant to the discussion we will sometimes drop the various annotations on these types; for example $\Gamma \vdash n : ch\langle U\rangle$ signifies that $ch_v\langle U\rangle$ for some visibility status v. Typing judgements take the form $\Gamma \vdash N$ and defined by the rules in Table 2. In these rules, we use an extended form of type environment, $\Sigma$, which, in addition to names, also maps *variables* to *stateless* types. Note that none of the rules depend on the status (dead or alive) of names in the environment. Also the visibility constraints are enforced indirectly, by virtue of the formation rules for valid types, given in Figure 1.

In this extended abstract we omit even the statement of the Subject Reduction and appropriate Type Safety result for our language.

Table 2. *Typing rules for typed DπLoc*

**Processes**

(t-out)
$$\frac{\begin{array}{c}\Sigma \vdash u : \mathrm{ch}\langle \tilde{\mathrm{U}}\rangle \\ \Sigma \vdash V : \tilde{\mathrm{U}} \\ \Sigma \vdash P\end{array}}{\Sigma \vdash u!\langle V\rangle.P}$$

(t-in-rep)
$$\frac{\begin{array}{c}\Sigma \vdash u : \mathrm{ch}\langle \tilde{\mathrm{U}}\rangle \\ \Sigma, X : \tilde{\mathrm{U}} \vdash P\end{array}}{\begin{array}{c}\Sigma \vdash u?(X).P \\ \Sigma \vdash *u?(X).P\end{array}}$$

(t-nw)
$$\frac{\Sigma, \tilde{n} : \tilde{\mathrm{T}} \vdash P}{\Sigma \vdash (\nu \tilde{n} : \tilde{\mathrm{T}})P}$$

(t-cond)
$$\frac{\begin{array}{c}\Sigma \vdash u : \mathrm{U}, \ v : \mathrm{U} \\ \Sigma \vdash P, \ Q\end{array}}{\Sigma \vdash \text{if } u = v \text{ then } P \text{ else } Q}$$

(t-fork)
$$\frac{\Sigma \vdash P, \ Q}{\Sigma \vdash P|Q}$$

(t-axiom)
$$\frac{}{\Sigma \vdash \mathbf{0}}$$

(t-go)
$$\frac{\begin{array}{c}\Sigma \vdash u : \mathrm{loc} \\ \Sigma \vdash P\end{array}}{\Sigma \vdash \mathrm{go}\ u.P}$$

(t-ping)
$$\frac{\begin{array}{c}\Sigma \vdash u : \mathrm{loc} \\ \Sigma \vdash P, \ Q\end{array}}{\Sigma \vdash \mathrm{ping}\ u.P \text{ else } Q}$$

**Systems**                                                    **Observers**

(t-rest)
$$\frac{\Gamma, \tilde{n} : \tilde{\mathrm{T}} \vdash N}{\Gamma \vdash (\nu \tilde{n} : \tilde{\mathrm{T}})N}$$

(t-par)
$$\frac{\Gamma \vdash N, \ M}{\Gamma \vdash N|M}$$

(t-proc)
$$\frac{\begin{array}{c}\Gamma \vdash l : \mathrm{loc} \\ \Gamma \vdash P\end{array}}{\Gamma \vdash l[\![P]\!]}$$

(t-obs)
$$\frac{\mathbf{pub}(\Gamma) \vdash O}{\Gamma \vdash_{\mathrm{obs}} O}$$

---

Table 3. *Reduction Rules for DπLoc*

Assuming $\Gamma \vdash l : \mathbf{alive}$

(r-comm)
$$\overline{\Gamma \rhd l[\![a!\langle V\rangle.P]\!] \mid l[\![a?(X).Q]\!] \quad \longrightarrow \quad \Gamma \rhd l[\![P]\!] \mid l[\![Q\{V\!/\!X\}]\!]}$$

(r-rep)
$$\overline{\Gamma \rhd l[\![*a?(X).P]\!] \quad \longrightarrow \quad \Gamma \rhd l[\![a?(X).(P \mid *a?(X).P)]\!]}$$

(r-fork)
$$\overline{\Gamma \rhd l[\![P|Q]\!] \quad \longrightarrow \quad \Gamma \rhd l[\![P]\!] \mid l[\![Q]\!]}$$

(r-eq)
$$\overline{\Gamma \rhd l[\![\text{if } u = u \text{ then } P \text{ else } Q]\!] \longrightarrow \Gamma \rhd l[\![P]\!]}$$

(r-neq)
$$\overline{\Gamma \rhd l[\![\text{if } u = v \text{ then } P \text{ else } Q]\!] \longrightarrow \Gamma \rhd l[\![Q]\!]} \; u \neq v$$

(r-go)
$$\overline{\Gamma \rhd l[\![\mathrm{go}\ k.P]\!] \longrightarrow \Gamma \rhd k[\![P]\!]} \; \Gamma \vdash k : \mathbf{alive}$$

(r-ngo)
$$\overline{\Gamma \rhd l[\![\mathrm{go}\ k.P]\!] \longrightarrow \Gamma \rhd k[\![\mathbf{0}]\!]} \; \Gamma \nvdash k : \mathbf{alive}$$

(r-ping)
$$\overline{\Gamma \rhd l[\![\mathrm{ping}\ k.P \text{ else } Q]\!] \longrightarrow \Gamma \rhd l[\![P]\!]} \; \Gamma \vdash k : \mathbf{alive}$$

(r-nping)
$$\overline{\Gamma \rhd l[\![\mathrm{ping}\ k.P \text{ else } Q]\!] \longrightarrow \Gamma \rhd l[\![Q]\!]} \; \Gamma \nvdash k : \mathbf{alive}$$

(r-new)
$$\overline{\Gamma \rhd l[\![(\nu n : \mathrm{T})P]\!] \longrightarrow \Gamma \rhd (\nu n : \mathrm{T})\ l[\![P]\!]}$$

(r-str)
$$\frac{\Gamma \rhd N' \equiv \Gamma \rhd N \quad \Gamma \rhd N \longrightarrow \Gamma' \rhd M \quad \Gamma' \rhd M \equiv \Gamma' \rhd M'}{\Gamma \rhd N' \longrightarrow \Gamma' \rhd M'}$$

(r-ctxt-rest)
$$\frac{\Gamma + n : \mathrm{T} \rhd N \quad \longrightarrow \quad \Gamma' + n : \mathrm{U} \rhd M}{\Gamma \rhd (\nu n : \mathrm{T})N \quad \longrightarrow \quad \Gamma' \rhd (\nu n : \mathrm{U})M}$$

(r-ctxt-par)
$$\frac{\Gamma \rhd N \quad \longrightarrow \quad \Gamma' \rhd N'}{\begin{array}{c}\Gamma \rhd N|M \quad \longrightarrow \quad \Gamma' \rhd N'|M \\ \Gamma \rhd M|N \quad \longrightarrow \quad \Gamma' \rhd M|N'\end{array}}$$

Table 4. *Structural Rules for DπLoc*

| | | |
|---|---|---|
| (s-comm) | $N\vert M \equiv M\vert N$ | |
| (s-assoc) | $(N\vert M)\vert M' \equiv N\vert(M\vert M')$ | |
| (s-unit) | $N\vert l[\![\mathbf{0}]\!] \equiv N$ | |
| (s-extr) | $(\nu\, n\!:\!\mathrm{T})(N\vert M) \equiv N\vert(\nu\, n\!:\!\mathrm{T})M$ | $n \notin \mathbf{fn}(N)$ |
| (s-flip) | $(\nu\, n\!:\!\mathrm{T})(\nu\, m\!:\!\mathrm{U})N \equiv (\nu\, m\!:\!\mathrm{U})(\nu\, n\!:\!\mathrm{T})N$ | |
| (s-inact) | $(\nu\, n\!:\!\mathrm{T})N \equiv N$ | $n \notin \mathbf{fn}(N)$ |

*Reduction Semantics:* We call pairs $\Gamma \triangleright N$ *configurations*, whenever $\Gamma \vdash N$. Reductions then take the form of a binary relation over configurations

$$\Gamma \triangleright N \;\;\longrightarrow\;\; \Gamma \triangleright N'$$

defined in terms of the reduction rules in Table 3, whereby systems reduce with respect to the status of the locations in $\Gamma$; we write $\Gamma \vdash l\!:\!\mathbf{alive}$ as a shorthand for $\Gamma \vdash l : \mathtt{loc}^{\mathrm{a}}$. So all reduction rules assume the location where the code is executing is alive. Moreover, (r-go), (r-ngo), (r-ping) and (r-nping) reduce according to the status of the remote location concerned. The reader is refered to [5] for more details; but note that here the status of locations is unchanged by reductions.

*Behavioural equivalence:* First note that the type system does indeed enforce the intuitive separation of concerns discussion in the Introduction. For example let $\Gamma_e$ denote the environment

$$\Gamma_e \;=\; l\!:\!\mathtt{loc}^{\mathrm{a}}_{\mathrm{p}},\; k_1\!:\!\mathtt{loc}^{\mathrm{a}}_{\mathrm{c}},\; k_2\!:\!\mathtt{loc}^{\mathrm{a}}_{\mathrm{c}},\; k_3\!:\!\mathtt{loc}^{\mathrm{a}}_{\mathrm{c}},\; req\!:\!\mathtt{ch}_{\mathrm{p}}\langle \mathrm{T}, \mathtt{ch}_{\mathrm{p}}\langle \mathrm{T}\rangle\rangle,\; a\!:\!\mathtt{ch}_{\mathrm{p}}\langle \mathrm{A}\rangle,\; ret\!:\!\mathtt{ch}_{\mathrm{p}}\langle \mathrm{T}\rangle$$

where T is an arbitrary public type; Then one can check

$$\Gamma_e \vdash \mathsf{server}_i$$

where $\mathsf{server}_i$ is defined in the Introduction, provided the locally declared channels *data* and $\mathsf{sync}$ are declared at the types $\mathtt{ch}\langle \mathrm{T}, \mathtt{ch}_{\mathrm{p}}\langle \mathrm{T}\rangle, \mathtt{loc}_{\mathrm{p}}\rangle$ and $\mathtt{ch}\langle \mathrm{T}\rangle$ respectively. Now consider

$$\mathsf{serverBad} \Leftarrow \mathsf{server}_1 \mid l[\![a!\langle k_1\rangle]\!]$$

which attempts to export a confined location $k_1$, which could subsequently could be tested for failure by a public observer. Once more one can check that $\Gamma_e \nvdash \mathsf{serverBad}$.

Intuitively an observer is any system which only uses public names. Formally let $\mathbf{pub}(\Gamma)$ be the environment obtained by omitting from $\Gamma$ any name not assigned a public type. Then $\mathbf{pub}(\Gamma) \vdash O$ ensures that $O$ can only use public names. For example consider

$$\mathsf{observer} \;\Leftarrow\; l[\![req!\langle v, ret\rangle]\!]$$
$$\mathsf{observerBad} \;\Leftarrow\; l[\![\mathsf{go}\, k_1.\mathsf{go}\, l.ok!\langle\rangle]\!]$$

Here one can check that $\mathbf{pub}(\Gamma_e) \vdash \mathsf{observer}$ and $\mathbf{pub}(\Gamma_e) \nvdash \mathsf{observerBad}$.

Our behavioural equivalence will in general relate arbitrary configurations; but we would expect equivalent configurations to have the same *public interface*, and be preserved by public observers.

**Definition 1 (p-Contextual).** *A relation over configurations is called p-Contextual if, whenever* $\Gamma \triangleright M \; \mathcal{R} \; \Gamma' \triangleright N$

- *(p-Interfaces:)* $\mathbf{pub}(\Gamma) = \mathbf{pub}(\Gamma')$
- *(Parallel:)* $\Gamma \triangleright M \mid O \; \mathcal{R} \; \Gamma' \triangleright N \mid O$ *and* $\Gamma \triangleright M \mid O \; \mathcal{R} \; \Gamma' \triangleright N \mid O$ *whenever* $\mathbf{pub}(\Gamma) \vdash O$
- *(Fresh extensions:)* $\Gamma, n : \mathsf{P} \triangleright M \; \mathcal{R} \; \Gamma', n : \mathsf{P} \triangleright N$ *whenever* $n$ *is fresh*

**Definition 2 (p-Barb).** $\Gamma \triangleright N \Downarrow^p_{a@l}$ *denotes a* p-observable barb *by configuration* $\Gamma \triangleright N$, *on channel a at location l, defined as:*

$\exists N'. \; \Gamma \triangleright N \longrightarrow^* \Gamma \triangleright N'$ *such that* $N' \equiv (\nu \, \tilde{n} : \tilde{\mathsf{T}}) M \mid l[\![a!\langle V \rangle.Q]\!]$ *where* $\Gamma \vdash l : \mathsf{loc}^{\mathsf{a}}_{\mathsf{p}}, a : \mathsf{ch}_{\mathsf{p}}\langle \tilde{\mathsf{P}} \rangle$

Using this concept, we can now modify the standard definition of *reduction barbed equivalence*, [6]:

**Definition 3 (Reduction barbed congruence).** *Let* $\cong$ *be the largest relation between configurations which is p-contextual, reduction-closed (see [6]) and preserves p-barbs.*

## 3 Defining Fault Tolerance

Our first notion of *n*-fault-tolerance, formalising the intuitive (1), is when the faulting context induces at most *n* location failures, prior to the execution of the system; of course these failures must only be induced on locations which are not public. Formally for any set of location names $\tilde{l}$ let $F^{\tilde{l}}_S$ be the function which maps any configuration $\Gamma \triangleright N$ to $\Gamma - \tilde{l} \triangleright N$, where $\Gamma - \tilde{l}$ is the environment obtained from $\Gamma$ by changing the status of every $l_i$ to *dead*. We say $F^{\tilde{l}}_S$ is a *valid static n-fault context* with respect to $\Gamma$, if the size of $\tilde{l}$ is *at most n*, and for every $l_i \in \tilde{l}$, $l_i$ is confined and alive ($\Gamma \vdash l_i : \mathsf{loc}^{\mathsf{a}}_{\mathsf{c}}$).

**Definition 4 (Static Fault Tolerance).** *A configuration* $\Gamma \triangleright N$ *is n-static fault tolerant if*

$$\Gamma \triangleright N \cong F^{\tilde{l}}_S(\Gamma) \triangleright N$$

*for every static n-fault context* $F^{\tilde{l}}_S$ *which is valid with respect to* $\Gamma$.

With this formal definition we can now examine the systems $\mathsf{server}_i$, using the $\Gamma_e$ defined above. We can easily check that $\Gamma \triangleright \mathsf{server}_1$ is not 1-fault tolerant, by considering the fault context $F^{k_1}_S$. Similarly we can show that $\Gamma_e \triangleright \mathsf{server}_2$ is not 2-fault tolerant, by considering $F^{k_1,k_2}_S$. But establishing positive results, for example that $\Gamma_e \triangleright \mathsf{server}_2$ is 1-fault tolerant, is difficult because the definition of $\cong$ quantifies over all valid observers. This point will be addressed in the next section, when we give a co-inductive characterisation of $\cong$.

Instead let us consider another manner of inducing faults. Let $l[\![\mathsf{kill}]\!]$ be a system which asynchronously kills a confined location $l$. Its operation is defined by the rule

(r-kill)

$$\frac{}{\Gamma \triangleright l[\![\mathsf{kill}]\!] \longrightarrow (\Gamma - l) \triangleright l[\![\mathbf{0}]\!]}$$

For any set of locations $\tilde{l}$ let $F^{\tilde{l}}_D$ denote the function which maps the system $M$ to $M \mid l_1[\![\mathsf{kill}]\!] \mid \ldots \mid l_n[\![\mathsf{kill}]\!]$. It is said to be a valid dynamic *n*-fault context with respect to $\Gamma$ if again the size of $\tilde{l}$ is at most *n* and $\Gamma \vdash l_i : \mathsf{loc}^{\mathsf{a}}_{\mathsf{c}}$, for every $l_i$ in $\tilde{l}$.

**Definition 5 (Dynamic Fault Tolerance).** *A configuration $\Gamma \triangleright N$ is n-dynamic fault tolerant if*

$$\Gamma \triangleright F_D^{\tilde{l}}(M) \cong \Gamma \triangleright M$$

*for every dynamic n-fault context which is valid with respect to $\Gamma$.*

*Example 2.* The system sPassive defined below uses two identical replicas of the distributed database at $k_1$ and $k_2$, but treats the replica at $k_1$ as *primary* replica and the one at $k_2$ as a *secondary* (backup) replica - once again $\mathtt{W} = ch\langle T, ch_p\langle T\rangle, loc_p\rangle$.

$$sPassive \Leftarrow (\nu\, data : \mathtt{W}) \left( l \left\|\begin{array}{l} serv?(x,y).ping\ k_1.go\ k_1.data!\langle x, y, l\rangle \\ \qquad\qquad\qquad\qquad else\ go\ k_2.data!\langle x, y, l\rangle \end{array}\right\| \begin{array}{l} \\ | \ k_1[\![data?(x,y,z).go\ z\ .y!\langle f(x)\rangle]\!] \\ | \ k_2[\![data?(x,y,z).go\ z\ .y!\langle f(x)\rangle]\!] \end{array}\right)$$

The coordinating interface at $l$ uses the ping construct to *detect failures* in the primary replica: if $k_1$ is alive, the request is sent to the primary replica and the secondary replica at $k_2$ is *not invoked*; if, on the other hand, the primary replica is dead, then the passive replica at $k_2$ is promoted to a primary replica and the request is sent to it. This implementation saves on *time redundancy* since, for any request, only one replica is invoked. Another passive replication server is sMonitor, defined as

$$sMonitor \Leftarrow (\nu\, data : \mathtt{W}) \left( \begin{array}{l} l \left\|\, serv?(x,y).(\nu\, sync : ch\langle\rangle) \left(\begin{array}{l} go\ k_1.data!\langle x, sync, l\rangle \\ | \ mntr\ k_1.go\ k_2.data!\langle x, sync, l\rangle \\ | \ sync?(z).y!\langle z\rangle \end{array}\right) \right\| \\ | \ k_1[\![data?(x,y,z).go\ z\ .y!\langle f(x)\rangle]\!] \\ | \ k_2[\![data?(x,y,z).go\ z\ .y!\langle f(x)\rangle]\!] \end{array}\right)$$

where again, $\mathtt{W} = ch\langle T, ch_p\langle T\rangle, loc_p\rangle$. It uses a *monitor* process for failure detection

$$mntr\ k.P \Leftarrow (\nu\, test : ch\langle\rangle)(\, test!\langle\rangle \,| \, * test?().ping\ k.\, test!\langle\rangle\ else\ P\,)$$

instead of a *single* ping test on the primary replica at $k_1$; mntr $k.P$ *repeatedly* tests the status of the monitored location $(k)$ and continues as $P$ when $k$ becomes dead. Similar to server$_{2..3}$, sMonitor synchronises multiple answers from replicas with channel *sync*.

Using the techniques of the next section, one can show that both $\Gamma_e \triangleright$ sPassive and $\Gamma \triangleright$ sMonitor, are 1-static fault tolerant, similar to server$_2$. However there is a difference between these two systems;if $k_1$ fails *after* sPassive tests for its status, then an answer will never reach $l$. Thus sPassive is *not* 1-dynamic fault tolerant; formally one can show $\Gamma_e \triangleright F_D^{k_1}(sPassive) \not\cong \Gamma_e \triangleright$ sPassive. But, as we will see in the next section, sMonitor can be shown to be 1-dynamic fault tolerant, just like server$_{2..3}$.

## 4 Proof Techniques for Fault Tolerance

We define a labelled transition system (lts) for D$\pi$Loc$_e$, which consists of a collection of actions over (closed) configurations, $\Gamma \triangleright N \xrightarrow{\mu} \Gamma' \triangleright N'$, where $\mu$ can be an internal action, $\tau$, a bound input, $(\tilde{n} : \tilde{T})l : a?(V)$ or bound output, $(\tilde{n} : \tilde{T})l : a!\langle V\rangle$. These actions

Table 5. *Operational Rules for Typed DπLoc*

---

Assuming $\Gamma \vdash l : \mathbf{alive}$

(l-in)
$$\frac{}{\Gamma \triangleright l[\![a?(X).P]\!] \xrightarrow{l:a?(V)} \Gamma \triangleright l[\![P\{V/X\}]\!]} \;\; \Gamma \vdash_{\mathsf{obs}} l, \; \Gamma \vdash a : \mathsf{ch}_p\langle \tilde{w}\rangle, \; V : \tilde{w}$$

(l-fork)
$$\frac{}{\Gamma \triangleright l[\![P|Q]\!] \xrightarrow{\tau} \Gamma \triangleright l[\![P]\!] \mid l[\![Q]\!]}$$

(l-out)
$$\frac{}{\Gamma \triangleright l[\![a!\langle V\rangle.P]\!] \xrightarrow{l:a!\langle V\rangle} \Gamma \triangleright l[\![P]\!]} \;\; \Gamma \vdash_{\mathsf{obs}} l, a$$

(l-in-rep)
$$\frac{}{\Gamma \triangleright l[\![*a?(X).P]\!] \xrightarrow{\tau} \Gamma \triangleright l[\![a?(X).(P \mid *a?(Y).P\{Y/X\})]\!]}$$

(l-eq)
$$\frac{}{\Gamma \triangleright l[\![\mathsf{if}\; u = u \;\mathsf{then}\; P \;\mathsf{else}\; Q]\!] \xrightarrow{\tau} \Gamma \triangleright l[\![P]\!]}$$

(l-neq)
$$\frac{}{\Gamma \triangleright l[\![\mathsf{if}\; u = v \;\mathsf{then}\; P \;\mathsf{else}\; Q]\!] \xrightarrow{\tau} \Gamma \triangleright l[\![Q]\!]} \;\; u \neq v$$

(l-new)
$$\frac{}{\Gamma \triangleright l[\![(\nu n : \mathsf{T})P]\!] \longrightarrow \Gamma \triangleright (\nu n : \mathsf{T})\, l[\![P]\!]}$$

(l-kill)
$$\frac{}{\Gamma \triangleright l[\![\mathsf{kill}]\!] \xrightarrow{\tau} (\Gamma - l) \triangleright l[\![\mathbf{0}]\!]}$$

(l-go)
$$\frac{}{\Gamma \triangleright l[\![\mathsf{go}\; k.P]\!] \xrightarrow{\tau} \Gamma \triangleright k[\![P]\!]} \;\; \Gamma \vdash k : \mathbf{alive}$$

(l-ngo)
$$\frac{}{\Gamma \triangleright l[\![\mathsf{go}\; k.P]\!] \xrightarrow{\tau} \Gamma \triangleright k[\![\mathbf{0}]\!]} \;\; \Gamma \nvdash k : \mathbf{alive}$$

(l-ping)
$$\frac{}{\Gamma \triangleright l[\![\mathsf{ping}\; k.P \;\mathsf{else}\; Q]\!] \xrightarrow{\tau} \Gamma \triangleright l[\![P]\!]} \;\; \Gamma \vdash k : \mathbf{alive}$$

(l-nping)
$$\frac{}{\Gamma \triangleright l[\![\mathsf{ping}\; k.P \;\mathsf{else}\; Q]\!] \xrightarrow{\tau} \Gamma \triangleright l[\![Q]\!]} \;\; \Gamma \nvdash k : \mathbf{alive}$$

(l-open)
$$\frac{\Gamma + n : \mathsf{T} \triangleright N \xrightarrow{(\tilde{n}:\tilde{\mathsf{T}})l:a!\langle V\rangle} \Gamma' \triangleright N'}{\Gamma \triangleright (\nu n : \mathsf{T})N \xrightarrow{(n:\mathsf{T},\tilde{n}:\tilde{\mathsf{T}})l:a!\langle V\rangle} \Gamma' \triangleright N'} \;\; l, a \neq n \in V$$

(l-weak)
$$\frac{\Gamma + n : \mathsf{T} \triangleright N \xrightarrow{(\tilde{n}:\tilde{\mathsf{T}})l:a?(V)} \Gamma' \triangleright N'}{\Gamma \triangleright N \xrightarrow{(n:\mathsf{T},\tilde{n}:\tilde{\mathsf{T}})l:a?(V)} \Gamma' \triangleright N'} \;\; l, a \neq n \in V$$

(l-rest)
$$\frac{\Gamma + n : \mathsf{T} \triangleright N \xrightarrow{\mu} \Gamma' + n : \mathsf{U} \triangleright N'}{\Gamma \triangleright (\nu n : \mathsf{T})N \xrightarrow{\mu} \Gamma' \triangleright (\nu n : \mathsf{U})N'} \;\; n \notin \mathbf{fn}(\mu)$$

(l-par-ctxt)
$$\frac{\Gamma \triangleright N \xrightarrow{\mu} \Gamma' \triangleright N'}{\Gamma \triangleright N \mid M \xrightarrow{\mu} \Gamma' \triangleright N' \mid M}$$
$$\frac{}{\Gamma \triangleright M \mid N \xrightarrow{\mu} \Gamma' \triangleright M \mid N'}$$

(l-par-comm)
$$\frac{\uparrow(\Gamma) \triangleright N \xrightarrow{(\tilde{n}:\uparrow(\tilde{\mathsf{T}}))l:a!\langle V\rangle} \Gamma' \triangleright N' \qquad \uparrow(\Gamma) \triangleright M \xrightarrow{(\tilde{n}:\uparrow(\tilde{\mathsf{T}}))l:a?(V)} \Gamma'' \triangleright M'}{\Gamma \triangleright N \mid M \xrightarrow{\tau} \Gamma \triangleright (\nu\tilde{n} : \tilde{\mathsf{T}})(N' \mid M')}$$
$$\frac{}{\Gamma \triangleright M \mid N \xrightarrow{\tau} \Gamma \triangleright (\nu\tilde{n} : \tilde{\mathsf{T}})(M' \mid N')}$$

---

are defined by transition rules given in Table 5, inspired by [7, 6, 5]. In accordance with Definition 2 (observable barbs) and Definitions 1 (valid observers), (l-in) and (l-out) restrict external communication to *public* channels at *public* locations ($\Gamma \vdash_{\mathsf{obs}} l, a$). Furthermore, in (l-in) we require that the type of the values inputted, $V$, match the object type of channel $a$; since $a$ is public and configurations are well-typed, this also implies that $V$ are public values defined in $\Gamma$. The restriction on output action, together with the assumption of well-typed configurations also means that, in (l-open), we only scope

**Table 6.** *β-Transition Rules for Typed DπLoc*

---

Assuming  $\Gamma \vdash l :$ **alive**

(b-in-rep)

$$\overline{\Gamma \triangleright l[\![*a?(X).P]\!] \stackrel{\tau}{\longmapsto}_\beta \Gamma \triangleright l[\![a?(X).(P| * a?(Y).P\{^Y\!/\!X\})]\!]}$$

(b-fork)

$$\overline{\Gamma \triangleright l[\![P|Q]\!] \stackrel{\tau}{\longmapsto}_\beta \Gamma \triangleright l[\![P]\!] \mid l[\![Q]\!]}$$

(b-eq)

$$\overline{\Gamma \triangleright l[\![\text{if } u = u \text{ then } P \text{ else } Q]\!] \stackrel{\tau}{\longmapsto}_\beta \Gamma \triangleright l[\![P]\!]}$$

(b-neq)

$$\overline{\Gamma \triangleright l[\![\text{if } u = v \text{ then } P \text{ else } Q]\!] \stackrel{\tau}{\longmapsto}_\beta \Gamma \triangleright l[\![Q]\!]} \; u \neq v$$

(b-ngo)

$$\overline{\Gamma \triangleright l[\![\text{go } k.P]\!] \stackrel{\tau}{\longmapsto}_\beta \Gamma \triangleright k[\![\mathbf{0}]\!]} \; \Gamma \nvdash k : \textbf{alive}$$

(b-nping)

$$\overline{\Gamma \triangleright l[\![\text{ping } k.P \text{ else } Q]\!] \stackrel{\tau}{\longmapsto}_\beta \Gamma \triangleright l[\![Q]\!]} \; \Gamma \nvdash k : \textbf{alive}$$

(b-new)

$$\overline{\Gamma \triangleright l[\![(\nu n : \text{T})P]\!] \stackrel{\tau}{\longmapsto}_\beta \Gamma \triangleright (\nu n : \text{T})l[\![P]\!]}$$

(b-rest)

$$\frac{\Gamma, n : \text{T} \triangleright N \stackrel{\tau}{\longmapsto}_\beta \Gamma', n : \mathbb{W} \triangleright N'}{\Gamma \triangleright (\nu n : \text{T})N \stackrel{\tau}{\longmapsto}_\beta \Gamma' \triangleright (\nu n : \mathbb{W})N'}$$

(b-par)

$$\frac{\Gamma \triangleright N \stackrel{\tau}{\longmapsto}_\beta \Gamma' \triangleright N'}{\Gamma \triangleright N|M \stackrel{\tau}{\longmapsto}_\beta \Gamma' \triangleright N'|M}$$
$$\Gamma \triangleright M|N \stackrel{\tau}{\longmapsto}_\beta \Gamma' \triangleright M|N'$$

---

extrude public values. Contrary to [5], the lts does not allow external killing of locations (through the label kill : *l*) since public locations are reliable and never fail. Finally, the transition rule for internal communication, (l-par-comm), uses an overloaded function $\uparrow ()$ for inferring input/output capabilities of the sub-systems: when applied to types, $\uparrow (\text{T})$ transforms all the type tags to public (p); when applied to environments, $\uparrow (\Gamma)$ changes all the types to public types in the same manner. All the remaining rules are a simplified version of the rules in [5].

**Definition 6 (Weak bisimulation equivalence).** *This is denoted as $\approx$, and is defined to be the largest typed relation over configurations such that if $\Gamma \triangleright M \approx \Gamma' \triangleright N$ then*

- *$\Gamma \triangleright M \stackrel{\mu}{\longrightarrow} \Gamma'' \triangleright M'$ implies $\Gamma' \triangleright N \stackrel{\hat{\mu}}{\Longrightarrow} \Gamma''' \triangleright N'$ such that $\Gamma'' \triangleright M' \approx \Gamma''' \triangleright N'$*
- *$\Gamma' \triangleright N \stackrel{\mu}{\longrightarrow} \Gamma'' \triangleright N'$ implies $\Gamma \triangleright M \stackrel{\hat{\mu}}{\Longrightarrow} \Gamma'' \triangleright M'$ such that $\Gamma'' \triangleright M' \approx \Gamma''' \triangleright N'$*

**Theorem 1 (Full Abstraction).** *For any DπLoc configurations $\Gamma \triangleright M$, $\Gamma' \triangleright N$:*

$$\Gamma \triangleright M \cong \Gamma' \triangleright N \quad \text{if and only if} \quad \Gamma \triangleright M \approx \Gamma' \triangleright N$$

Theorem 1 allows us to prove *positive* fault tolerance results by giving a bisimulation for every reduction barbed congruent pair required by Definitions 4 and 5. We next develop up-to bisimulation techniques that can relieve some of the burden of exhibiting the required bisimulations. We identify a number of $\tau$ actions, which we refer to as $\beta$-actions or $\beta$-moves, inspired by the work in [3]. These are denoted as $\Gamma \triangleright N \stackrel{\tau}{\longmapsto}_\beta \Gamma' \triangleright N$ and are defined in Table 6. With these $\beta$-moves we develop up-to bisimulation techniques, by showing that our witness bisimulations can abstract away from matching configurations that denote $\beta$-moves. Our details are more complicated than in [3] because we deal with failure: apart from local rules ((b-eq) and (b-fork)) and context rules ((b-rest) and (b-par)), Table 6 includes rules dealing with failed locations such as (b-ngo)

Table 7. *β-Equivalence Rules for Typed DπLoc*

| | | |
|---|---|---|
| (bs-comm) | $\Gamma \models N|M \equiv_f M|N$ | |
| (bs-assoc) | $\Gamma \models (N|M)|M' \equiv_f N|(M|M')$ | |
| (bs-unit) | $\Gamma \models N|l[\![\mathbf{0}]\!] \equiv_f N$ | |
| (bs-extr) | $\Gamma \models (\nu n\!:\!\mathrm{T})(N|M) \equiv_f N|(\nu n\!:\!\mathrm{T})M$ | $n \notin \mathbf{fn}(N)$ |
| (bs-flip) | $\Gamma \models (\nu n\!:\!\mathrm{T})(\nu m\!:\!\mathrm{U})N \equiv_f (\nu m\!:\!\mathrm{U})(\nu n\!:\!\mathrm{T})N$ | |
| (bs-inact) | $\Gamma \models (\nu n\!:\!\mathrm{T})N \equiv_f N$ | $n \notin \mathbf{fn}(N)$ |
| (bs-dead) | $\Gamma \models l[\![P]\!] \equiv_f l[\![Q]\!]$ | $\Gamma \nvdash l : \mathbf{alive}$ |

and (b-nping). To obtain the required results for $\beta$-moves with failure, we define a new structural equivalence ranging over *configurations*, denoted as $\equiv_f$ and defined by the rules in Table 7, which takes into consideration *location status* as well. This enables us to obtain confluence for $\beta$-moves with respect to actions that change the status of locations. The only rule worth highlighting is (bs-dead), which allows us to ignore dead code.

**Lemma 1 (Confluence of $\beta$-moves).** $\overset{\tau}{\longmapsto}_\beta$ *observes the diamond property:*

$$\Gamma \rhd N \overset{\tau}{\underset{\beta}{\longmapsto}} \Gamma \rhd M \quad \text{implies} \quad \Gamma \rhd N \overset{\tau}{\underset{\beta}{\longmapsto}} \Gamma \rhd M \quad \text{or } \mu = \tau \text{ and } \Gamma \rhd M = \Gamma' \rhd N'$$

$$\mu \downarrow \qquad\qquad\qquad\qquad \mu \downarrow \qquad \mu \downarrow$$

$$\Gamma' \rhd N' \qquad \Gamma' \rhd M' \qquad \Gamma' \rhd N' \overset{\tau}{\underset{\beta}{\Longrightarrow}} \equiv_f \Gamma' \rhd M'$$

*Proof.* The proof proceeds by case analysis of the different types of $\mu$ and then by induction on the derivation of the $\beta$-move.

**Proposition 1.** *Suppose $\Gamma \rhd N \Longrightarrow_\beta \Gamma' \rhd M$. Then $\Gamma \rhd N \approx \Gamma' \rhd M$.*

*Proof.* We prove the above statement by defining $\mathcal{R} = \left\{ \Gamma \rhd N , \Gamma' \rhd M \mid \Gamma \rhd N \Longrightarrow_\beta \Gamma' \rhd M \right\}$ and showing that $\mathcal{R}$ is a bisimulation, which follows as a consequence of Lemma 1.

**Definition 7 (Bisimulation up-to $\beta$-moves).** *Bisimulation up-to $\beta$-moves, denoted as $\approx_\beta$, is the largest typed relation between configurations such that $\Gamma_1 \rhd M_1 \approx_\beta \Gamma_2 \rhd M_2$ and*

- $\Gamma_1 \rhd M_1 \overset{\mu}{\longrightarrow} \Gamma'_1 \rhd M'_1$ *implies* $\Gamma_2 \rhd M_2 \overset{\hat{\mu}}{\Longrightarrow} \Gamma'_2 \rhd M'_2$ *such that* $\Gamma'_1 \rhd M'_1 \mathcal{A}_l \circ \approx_\beta \circ \approx \Gamma'_2 \rhd M'_2$
- $\Gamma_2 \rhd M_2 \overset{\mu}{\longrightarrow} \Gamma'_2 \rhd M'_2$ *implies* $\Gamma_1 \rhd M_1 \overset{\hat{\mu}}{\Longrightarrow} \Gamma'_1 \rhd M'_1$ *such that* $\Gamma'_2 \rhd M'_2 \mathcal{A}_l \circ \approx_\beta \circ \approx \Gamma'_1 \rhd M'_1$

*where $\mathcal{A}_l$ is the relation $\Longrightarrow_\beta \circ \equiv$.*

Proposition 1 provides us with a powerful method for approximating bisimulations. In the approximate bisimulation $\approx_\beta$, an action $\Gamma_1 \rhd M_1 \overset{\mu}{\longrightarrow} \Gamma'_1 \rhd M'_1$ can be matched by a $\beta$-derivative of $\Gamma'_1 \rhd M'_1$, that is $\Gamma'_1 \rhd M'_1 \Longrightarrow_\beta \Gamma'_1 \rhd M''_1$, and a weak matching action $\Gamma_2 \rhd M_2 \overset{\hat{\mu}}{\Longrightarrow} \Gamma'_2 \rhd M'_2$ such that, up to structural equivalence on the one side and up-to bisimilarity on the other, the pairs $\Gamma'_1 \rhd M''_1$ and $\Gamma'_2 \rhd M'_2$ are once more related. Intuitively then, in any relation satisfying $\approx_\beta$, a configuration can represent all the configurations to which it can evolve using $\beta$-moves. We justify the use of $\approx_\beta$ by proving Proposition 2.

**Proposition 2 (Inclusion of bisimulation up-to $\beta$-moves).** *If $\Gamma_1 \triangleright M_1 \approx_\beta \Gamma_2 \triangleright M_2$ then $\Gamma_1 \triangleright M_1 \approx \Gamma_2 \triangleright M_2$*

*Proof.* We prove the above proposition by defining the relation $\mathcal{R}$ as

$$\mathcal{R} = \left\{ \Gamma_1 \triangleright M_1 \,, \Gamma_2 \triangleright M_2 \quad \middle| \Gamma_1 \triangleright M_1 \approx \circ \approx_\beta \circ \approx \Gamma_2 \triangleright M_2 \right\}$$

and show that $\mathcal{R} \subseteq \approx$. The required result can then be extracted from this result by considering the special cases where the $\approx$ on either side are the identity relations.

*Example 3.* We are now in a position to prove positive fault tolerance result. For instance to show that $\Gamma \triangleright \mathsf{sPassive}$ is 1-static fault tolerant we just need to provide 3 witness bisimulations up-to $\beta$-moves to prove

$$\prod_{i=1}^{3} \quad \Gamma \triangleright \mathsf{sPassive} \quad \cong \quad (\Gamma - k_i) \triangleright \mathsf{sPassive}$$

We here give the witness relation for the most involving case (where $i = 1$), and leave the simpler relations for the interested reader. Thus, the witness relation is $\mathcal{R}$ defined as

$$\mathcal{R} \stackrel{\text{def}}{=} \{\langle \Gamma \triangleright \mathsf{sPassive}, \Gamma - k_1 \triangleright \mathsf{sPassive}\rangle\} \,\cup\, \left( \bigcup_{u,v \in \mathrm{N_{AMES}}} \mathcal{R}'(u,v) \right)$$

$$\mathcal{R}'(u,v) \stackrel{\text{def}}{=} \begin{cases} \Gamma \triangleright (\nu d)l[\![\mathsf{Png}(u,v)]\!] \mid \mathsf{R}_1 \mid \mathsf{R}_2 & ,\Gamma - k_1 \triangleright (\nu d)l[\![\mathsf{Png}(u,v)]\!] \mid \mathsf{R}_1 \mid \mathsf{R}_2 \\ \Gamma \triangleright (\nu d)l[\![\mathsf{Q}_1(u,v)]\!] \mid \mathsf{R}_1 \mid \mathsf{R}_2 & ,\Gamma - k_1 \triangleright (\nu d)l[\![\mathsf{Q}_2(u,v)]\!] \mid \mathsf{R}_1 \mid \mathsf{R}_2 \\ \Gamma \triangleright (\nu d)k_1[\![d!\langle u,v,l\rangle]\!] \mid \mathsf{R}_1 \mid \mathsf{R}_2 & ,\Gamma - k_1 \triangleright (\nu d)k_2[\![d!\langle u,v,l\rangle]\!] \mid \mathsf{R}_1 \mid \mathsf{R}_2 \\ \Gamma \triangleright (\nu d)k_1[\![\mathsf{go}\, l\,.v!\langle f(u)\rangle]\!] \mid \mathsf{R}_2 & ,\Gamma - k_1 \triangleright (\nu d)\mathsf{R}_1 \mid k_2[\![\mathsf{go}\, l\,.v!\langle f(u)\rangle]\!] \\ \Gamma \triangleright (\nu d)l[\![v!\langle f(u)\rangle]\!] \mid \mathsf{R}_2 & ,\Gamma - k_1 \triangleright (\nu d)\mathsf{R}_1 \mid l[\![v!\langle f(u)\rangle]\!] \\ \Gamma \triangleright (\nu d)\mathsf{R}_2 & ,\Gamma - k_1 \triangleright (\nu d)\mathsf{R}_1 \end{cases}$$

where $d$ stands for *data* and
$\mathsf{Png}(x,y) \Leftarrow \mathsf{ping}\, k_1.\mathsf{Q}_1(x,y)\,\mathsf{else}\,\mathsf{Q}_2(x,y)$
$\quad \mathsf{Q}_i(x,y) \Leftarrow \mathsf{go}\, k_i.d!\langle x,y,l\rangle$
$\qquad\quad \mathsf{R}_i \Leftarrow k_i[\![d?(x,y,z).\mathsf{go}\, z\,.y!\langle f(x)\rangle]\!]$

## 5 Generic Techniques for Dynamic Fault Tolerance

Despite the fault tolerance proof techniques developed in Section 4, proving positive fault tolerance results entails a lot of unnecessary repeated work because Definition 4 and Definition 5 quantify over all valid fault contexts: to prove that $\mathsf{server}_3$ is 2-dynamic fault tolerant, we need to provide 6 relations, one for every different case in

$$\prod_{i \neq j = 1}^{3} \quad \Gamma \triangleright \mathsf{server}_3 \quad \cong \quad \Gamma \triangleright \mathsf{server}_3 | k_i[\![\mathsf{kill}]\!] | k_j[\![\mathsf{kill}]\!]$$

A closer inspection of the required relations reveals that there is a lot of overlap between them: these overlapping states would be automatically circumvented if we require a single relation that is somewhat the merging of all of these separate relations. Hence we

reformulate our fault tolerance definition for dynamic fault tolerance (the most demanding), to reflect such a merging of relations; a similar definition for the static case should not be more difficult to construct. The new definition is based on the actions given earlier in Section 4 together with a new action, fail, defined as

$$(\text{l-fail})$$
$$\frac{}{\Gamma \triangleright N \xrightarrow{\text{fail}} (\Gamma - l) \triangleright N} \; \Gamma \vdash l\!:\!\texttt{loc}_{\texttt{c}}^{\texttt{a}}$$

permitting external killing of confined locations. Intuitively, this action allow us to *count* the number of failures, but prohibits us from determining which specific location failed.[3] The asymmetric relation $\preceq_D^n$, defined below, is parameterised with an integer $n$, denoting the number of confined locations that can still be killed on the right hand side: the additional third clause states that a fail move on the right hand side may be matched by a weak $\tau$-move on the left hand side and the two residuals need to be related in $\preceq_D^{n-1}$.

**Definition 8 (Dynamic Fault Tolerance Simulation).** *Dynamic n-fault tolerant simulation, denoted $\preceq_D^n$, is the largest asymmetric typed relation over configurations such that whenever $\Gamma_1 \triangleright M_1 \preceq_D^n \Gamma_2 \triangleright M_2$,*

- $\Gamma_1 \triangleright M_1 \xrightarrow{\gamma} \Gamma_1' \triangleright M_1'$ *implies* $\Gamma_2 \triangleright M_2 \stackrel{\widehat{\gamma}}{\Longrightarrow} \Gamma_2' \triangleright M_2'$ *such that* $\Gamma_1' \triangleright M_1' \preceq_D^n \Gamma_2' \triangleright M_2'$
- $\Gamma_2 \triangleright M_2 \xrightarrow{\gamma} \Gamma_2' \triangleright M_2'$ *implies* $\Gamma_1 \triangleright M_1 \stackrel{\widehat{\gamma}}{\Longrightarrow} \Gamma_1' \triangleright M_1'$ *such that* $\Gamma_1' \triangleright N_1' \preceq_D^n \Gamma_2' \triangleright M_2'$
- *if* $n > 0$, $\Gamma_2 \triangleright M_2 \xrightarrow{\text{fail}} \Gamma_2' \triangleright M_2'$ *implies* $\Gamma_1 \triangleright M_1 \Longrightarrow \Gamma_1' \triangleright M_1'$ *such that* $\Gamma_1' \triangleright M_1' \preceq_D^{n-1} \Gamma_2' \triangleright M_2'$

Before we can use Definition 8 to prove dynamic fault tolerance, we need to show that the new definition is sound with respect to Definition 5.

**Proposition 3 (Soundness of $\preceq_D^n$).** *If $\Gamma \models M_1 \; \preceq_D^n \; M_2$ then for any dynamic n-fault context $F_D^{\tilde{l}}$ that is valid with respect to $\Gamma$ we have $\Gamma \models M_1 \cong F_D^{\tilde{l}}(M_2)$*

*Proof.* Let $\mathcal{R}_n$ be a relation parameterised by a number $n$ and defined as

$$\mathcal{R}_n \stackrel{\text{def}}{=} \left\{ \Gamma_1 \triangleright M_1 \,, \Gamma_2 \triangleright M_2 \,|\, F_D^i \quad \Big| \Gamma_1 \triangleright M_1 \preceq_D^i \Gamma_2 \triangleright M_2, \; \textstyle\prod_{j=1}^2 \Gamma_j \vdash F_D^i \text{ and } 0 \leq i \leq n \right\}$$

By showing $\mathcal{R}_n \subseteq\approx$ we prove that $\preceq_D^n$ is sound with respect to $n$-dynamic fault tolerance

It would be ideal if we could reuse up-to techniques and give relations satisfying $\preceq_D^n$ that abstract away from $\beta$-moves. Similar to Section 4, we define a fault tolerance simulation up-to $\beta$-moves and show that this is sound with respect to $\preceq_D^n$. This definition uses a weak bisimulation (Definition 6) that ranges over $\alpha$ actions, that is $\mu$ *and* the new action fail. We refer to this bisimulation as a *counting* bisimulation over configurations, denoted as $\approx_{cnt}$, because it allows us to count failing confined locations on each side and match subsequent observable behaviour.

**Definition 9 (Fault Tolerant Simulation up-to $\beta$-moves).** *An n-fault tolerant simulation up-to $\beta$-moves, denoted as $\preceq_\beta^n$, is the largest typed relation $\mathcal{R}$ between configurations parameterised by the number n, such that whenever we have $\Gamma_1 \triangleright M_1 \preceq_\beta^n \Gamma_2 \triangleright M_2$*

---

[3] This point differs from [5], where labels for external killing carried the location name, kill:l.

- $\Gamma_1 \triangleright M_1 \xrightarrow{\mu} \Gamma_1' \triangleright M_1'$ *implies* $\Gamma_2 \triangleright M_2 \xLongrightarrow{\hat{\mu}} \Gamma_2' \triangleright M_2'$ *such that* $\Gamma_1' \triangleright M_1' \mathcal{A}_l \circ \preceq_\beta^n \circ \approx_{cnt} \Gamma_2' \triangleright M_2'$

- $\Gamma_2 \triangleright M_2 \xrightarrow{\mu} \Gamma_2' \triangleright M_2'$ *implies* $\Gamma_1 \triangleright M_1 \xLongrightarrow{\hat{\mu}} \Gamma_1' \triangleright M_1'$ *such that* $\Gamma_2' \triangleright M_2' \mathcal{A}_l \circ \preceq_\beta^n \circ \approx \Gamma_1' \triangleright M_1'$

- *If* $n > 0$ *then* $\Gamma_2 \triangleright M_2 \xrightarrow{fail} \Gamma_2' \triangleright M_2'$ *implies* $\Gamma_1 \triangleright M_1 \Longrightarrow \Gamma_1' \triangleright M_1'$ *such that*
$\Gamma_2' \triangleright M_2' \preceq_\beta^{n-1} \circ \approx \Gamma_1' \triangleright M_1'$

*where* $\mathcal{A}_l$ *is the relation* $\Longrightarrow_\beta \circ \equiv$. *We highlight the use of* $\approx_{cnt}$ *for matching configurations in the first clause.*

The work required to show that $\preceq_\beta^n$ is sound with respect to $\preceq_D^n$ is similar to earlier up-to $\beta$-moves work discussed in Section 4: we have to show that $\beta$-move confluence (similar to Lemma 1) is also preserved for the new action fail; we also have to show that after a $\beta$-move, the redex and reduct configurations are counting-bisimilar (similar to Proposition 1). Finally we prove the following proposition

**Proposition 4 (Inclusion of fault tolerant simulation up-to $\beta$-moves).**
*If* $\Gamma_1 \triangleright M_1 \preceq_\beta^n \Gamma_2 \triangleright M_2$ *then* $\Gamma_1 \triangleright M_1 \preceq_D^n \Gamma_2 \triangleright M_2$

*Proof.* We prove the above proposition by defining the relation $\mathcal{R}_n$ as

$$\mathcal{R}_n = \left\{ \Gamma_1 \triangleright M_1, \Gamma_2 \triangleright M_2 \quad \middle| \Gamma_1 \triangleright M_1 \approx \circ \preceq_\beta^i \circ \approx_{cnt} \Gamma_2 \triangleright M_2 \text{ and } 0 \le i \le n \right\}$$

and show that $\mathcal{R}_n \subseteq \preceq_D^n$. The required result can then be extracted from this result by considering the special cases where $\approx$ and $\approx_{cnt}$ on either side are the identity relations.

*Example 4.* The results of Proposition 3 and Proposition 4 allow us to prove that the configuration $\Gamma \triangleright \mathsf{server}_2$ is 1-dynamically fault tolerant by providing a *single* witness fault tolerance simulation up-to $\beta$-moves showing that $\Gamma \triangleright \mathsf{server}_2 \preceq_\beta^1 \Gamma \triangleright \mathsf{server}_2$ Due to lack of space, we relegate the presentation of this relation to the full paper [4].

## 6 Conclusions and Related Work

We adopted a subset of [5] and developed a theory for system fault tolerance in the presence of fail-stop node failure. We formalised two definitions for fault tolerance based on the well studied concept of observational equivalence. Subsequently, we developed various sound proof techniques with respect to these definitions.

*Future Work* The immediate next step is to apply the theory to a wider spectrum of examples, namely using replicas with state and fault tolerance techniques such as lazy replication: we postulate that the existing theory should suffice. Another avenue worth considering is extending the theory to deal with link failure and the interplay between node and link failure [5]. In the long run, we plan to develop of a compositional theory of fault tolerance, enabling the construction of fault tolerant systems from smaller component sub-systems. For both cases, this paper should provide a good starting point.

*Related Work* To the best of our knowledge, Prasad's thesis [9] is the closest work to ours, addressing fault tolerance for process calculi. Even though similar concepts such as redundancy (called "duplication") and failure-free execution are identified, the setting and development of Prasad differs considerably form ours. In essence, three new operators ("displace", "audit" and "checkpoint") are introduced in a CCS variant; equational laws for terms using these operators are then developed so that algebraic manipulation can be used to show that terms in this calculus are, in some sense, fault tolerant with respect to their specification.

## References

1. Roberto M. Amadio and Sanjiva Prasad. Localities and failures. *FSTTCS: Foundations of Software Technology and Theoretical Computer Science*, 14, 1994.
2. Flavin Christian. Understanding fault tolerant distributed systems. *Communications of the ACM*, 34(2):56–78, February 1991.
3. Alberto Ciaffaglione, Matthew Hennessy, and Julian Rathke. Proof methodologies for behavioural equivalence in D$\pi$. Technical Report 03/2005, University of Sussex, 2005.
4. Adrian Francalanza and Matthew Hennessy. A theory for observational fault tolerance. `www.cs.um.edu.mt/˜ afran/`.
5. Adrian Francalanza and Matthew Hennessy. A theory of system behaviour in the presence of node and link failures. In *CONCUR*, volume 3653 of *Lecture Notes in Computer Science*, pages 368–382. Springer, 2005.
6. Matthew Hennessy, Massimo Merro, and Julian Rathke. Towards a behavioural theory of access and mobility control in distributed systems. *Theoretical Computer Science*, 322:615–669, 2004.
7. Matthew Hennessy and Julian Rathke. Typed behavioural equivalences for processes in the presence of subtyping. *Mathematical Structures in Computer Science*, 14:651–684, 2004.
8. Matthew Hennessy and James Riely. Resource access control in systems of mobile agents. *Information and Computation*, 173:82–120, 2002.
9. K. V. S. Prasad. *Combinators and Bisimulation Proofs for Restartable Systems*. PhD thesis, Department of Computer Science, University of Edinburgh, December 1987.
10. James Riely and Matthew Hennessy. Distributed processes and location failures. *Theoretical Computer Science*, 226:693–735, 2001.
11. Davide Sangiorgi and David Walker. *The $\pi$-calculus*. Cambridge University Press, 2001.
12. Richard D. Schlichting and Fred B. Schneider. Fail-stop processors: An approach to designing fault-tolerant computing systems. *Computer Systems*, 1(3):222–238, 1983.
13. Paulo Verissimo and Luis Rodrigues. *Distributed Systems for System Architects*. Kluwer Academic Publishers, 2001.