

# NeXeme: A Distributed Scheme Based on Nexus

Luc Moreau\*, David De Roure\*, and Ian Foster\*\*

**Abstract.** The remote service request, a form of remote procedure call, and the global pointer, a global naming mechanism, are two features at the heart of Nexus, a library for building distributed systems. NeXeme is an extension of Scheme that fully integrates both concepts in a mostly-functional framework, hence providing an expressive language for distributed computing. This paper presents a semantics for this Scheme extension, and also describes a NeXeme implementation, including its distributed garbage collector.

## 1 Introduction

Scheme [20] is a mostly-functional language, i.e., it is a fully functional language that also supports imperative notions such as assignments and continuations, for efficiency and expressivity reasons. We believe that a distributed extension of such a language requires a mechanism to invoke functions remotely, so that distribution becomes part of the most fundamental operation of the language. Such an approach is also adopted by languages such as Obliq [2], Java + RMI [10, 22], and Compositional C++ [4], in which methods can be invoked remotely.

Nexus [7], a library for building distributed systems, has two salient features: a *remote service request* is a form of remote procedure call [1], and *global pointers* provide for global naming in a distributed environment. By offering a functionality close to remote function invocation, Nexus is a suitable building block for our distributed language. Furthermore, when designing a distributed version of Scheme, our concerns were portability and potential use of high-performance hardware or protocols (e.g. supercomputers, ATM, UDP). Nexus also addresses these concerns as it runs on a variety of platforms and protocols.

NeXeme integrates the Nexus approach, with its remote service requests and global pointers, into a mostly functional language. The result is a novel distributed programming language that offers expressivity, development ease, and automatic memory management (via a distributed garbage collector). NeXeme provides powerful abstractions for controlling distribution while remaining computationally efficient. We believe that NeXeme is an excellent medium for implementing other forms of parallelism such as communication channels [21, 9] and futures [12, 14]. It is also an ideal platform for developing distributed symbolic applications, based for example on distributed mobile agents.

In this paper, we formalise the concepts of remote service request and global pointers. To this end, in Section 3, we present a formal semantics for a simplified version of NeXeme, called *Idealised NeXeme*. The semantics is operational as it defines a mechanical way of evaluating NeXeme programs on an abstract machine. In Section 4, we describe the NeXeme implementation. We compare and discuss our approach with related work in Section 5. More information on NeXeme is available at the following URL [15].

## 2 The Nexus Architecture

Nexus [7] is structured in terms of five basic abstractions: nodes, contexts, threads, global pointers, and remote service requests. A computation executes on a set of *nodes*

\* This research was supported in part by EPSRC grant GR/K30773. Authors' address: Department of Electronics and Computer Science, University of Southampton, Southampton SO17 1BJ, United Kingdom. E-mail: (L.Moreau,dder)@ecs.soton.ac.uk.

\*\* Author's address: Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, IL 60439. E-mail: foster@mcs.anl.gov.

and consists of a set of *threads*, each executing in an address space called a *context*. (For the purposes of this article, it suffices to assume that a context is equivalent to a process.) An individual thread executes a sequential program, which may read and write data shared with other threads executing in the same context.

The *global pointer* (GP) provides a global name space for objects, while the *remote service request* (RSR) is used to initiate communication and invoke remote computation. A GP represents a communication endpoint: that is, it specifies a destination to which a communication operation can be directed by an RSR. GPs can be created dynamically; once created, a GP can be communicated between nodes by including it in an RSR. A GP can be thought of as a capability granting rights to operate on the associated endpoint.

Practically, an RSR is specified by providing a global pointer, a handler identifier, and a data buffer, in which data are serialised. Issuing an RSR causes the data buffer to be transferred to the context designated by the global pointer, after which the routine specified by the handler is executed, potentially in a new thread of control. Both the data buffer and pointed specific data are available to the RSR handler.

The remote service request mechanism allows point-to-point communication, remote memory access, and streaming protocols to be supported within a single framework.

### 3 The Essence of Remote Service Requests

In this section, we present an operational semantics for remote service requests. For the sake of clarity and simplicity, we consider a subset of NeXeme, called idealised NeXeme. Table 1 summarises its different data structures, their representation, their constructors, and the operations permitted on them. (Programming examples may be found in [15].) The semantics was designed also to model non-functional languages like C; for example, full closures and first-class continuations are present in the semantics, but are not required in the actual implementation.

name	representation	constructors	operations
pair	$(\text{cons } V_1 \ V_2)$	cons	car, cdr
box	$\langle \text{bx } \alpha, s \rangle$	makeref	deref, setref
lock	$\langle \text{lk } \alpha, s \rangle$	make-lock	lock, unlock
global pointer	$\langle \text{gp } \alpha, s \rangle$	make-gp, remote	remote service request
closure	$\langle \text{cl } \lambda \vec{x}. M, \rho \rangle$	$\lambda \vec{x}. M$	apply
continuation	$\langle \text{co } \kappa \rangle$	callcc	apply

Table 1. Data Structures

The state space of the semantics appears in Figure 1. Distributed computations proceed inside a *world* composed of several *sites*<sup>3</sup>. Each site is identified by a unique name  $s$  and is composed of a set  $T$  of *tasks* sharing access to a *store*  $\theta$ . A store associates locations with values. Sites communicate by exchanging *requests*; each site contains both an output queue  $O$  of requests to be sent, and an input queue  $I$  of requests to be handled. When arriving at a site, requests are dealt with by a handler contained in  $H$ , a handler table associating names with closures.

A task is an abstraction of a Nexus thread composed of a *computational state* and a name. A computational state is a configuration of the CEK abstract machine [5] designed to evaluate the sequential subset of the language. Two states are permitted:  $\text{Ev}(M, \rho, \kappa)$  represents the evaluation of a term  $M$ , in the environment  $\rho$ , with a continuation  $\kappa$ ;  $\text{Ret}(V, \kappa)$  models the return of a value  $V$  to a continuation  $\kappa$ . In terms of implementation, the term being evaluated is the program counter, the continuation is the control stack, and the environment is the sequence of frames containing bindings.

<sup>3</sup> In the semantics, a site corresponds to the notion of *context* in Nexus; we adopt a different name because contexts are usually given a different meaning in semantics [5].

---

$\mathcal{W}$	$::= \{m_1, \dots, m_n\}$	(World)	Free Task Name:
$m \in \mathcal{M}$	$::= \langle T, \theta, s, H, I, O \rangle$	(Site)	$FTN(T) = \{\tau \mid \langle C, \tau \rangle \in T\}$
$t \in Task$	$::= \langle C, \tau \rangle$	(Task)	Free Site Name:
$\theta \in Store$	$::= \{(\alpha_1 V_1) \dots (\alpha_n V_n)\}$	(Store)	$FSN(\mathcal{W}) = \{s \mid \langle T, \theta, s, H, I, O \rangle \in \mathcal{W}\}$
$s \in \mathcal{S}$	$= \{s_1, s_2, \dots\}$	(Site Name)	
$\tau \in \mathcal{T}$	$= \{\tau_1, \tau_2, \dots\}$	(Task Name)	
$C \in CoSt$	$::= Ev\langle M, \rho, \kappa \rangle$ (Computational State)		$\delta(\text{cons}, V_1, V_2) = (\text{cons } V_1 V_2)$
	$  Ret\langle V, \kappa \rangle$		$\delta(\text{car}, (\text{cons } V_1 V_2)) = V_1$
$\rho \in Env$	$::= \{(x_1 V_1) \dots (x_n V_n)\}$ (Environment)		$\delta(\text{cdr}, (\text{cons } V_1 V_2)) = V_2$
$I$	$= \{R_1, \dots, R_n\}$ (Input Queue)		$\delta(+, [n], [m]) = [n + m]$
$O$	$= \{R_1, \dots, R_n\}$ (Output Queue)		
$T$	$= \{t_1, \dots, t_n\}$ (Tasks Set)		
$H \in Hndl$	$::= \{(str_1 f_1) \dots (str_n f_n)\}$ (Handlers)		
<hr/>			
$M \in \Lambda$	$::= V_s \mid (M M \dots) \mid (\text{if } M M M)$		(Term)
$V_s \in SValue$	$::= x \mid (\lambda \vec{x}. M) \mid str \mid c$		(Syntactic Value)
$V \in RValue$	$::= f \mid str \mid c \mid p \mid b \mid l \mid gp \mid k \mid \tau$		(Runtime Value)
$x \in Vars$	$= \{x, y, z, \dots\}$		(User Variable)
$\vec{x} \in VVec$	$::= x \mid x \vec{x} \mid .x \mid .0$		(Var Vector)
$f \in Clo$	$::= \langle cl \lambda \vec{x}. M, \rho \rangle$		(Closure)
$str \in String$			(String)
$c \in Const$	$::= c_b \mid c_f$		(Constant)
$c_b \in BConst$	$= \{\text{true}, \text{false}, \text{nil}, 0, 1, \dots, \text{void}\}$		(Basic Constant)
$c_f \in FConst$	$= \{\text{cons}, \text{car}, \text{cdr}, +, \text{makeref}, \text{deref},$ $\text{setref}, \text{callcc}, \text{make-gp}, \text{remote}, \text{rsr}, \text{fork}, \text{define-handler}, \text{shutdown}\}$		(Functional Constant)
$p \in Pair$	$::= (\text{cons } V V)$		(Pair)
$b \in Box$	$::= \langle bx \alpha, s \rangle$		(Box)
$l \in Locks$	$::= \langle lk \alpha, s \rangle$		(Lock)
$gp \in GloP$	$::= \langle gp s, \alpha \rangle$		(Global Pointer)
$k \in Cont$	$::= \langle co \kappa \rangle$		(Continuation)
$\kappa \in CCode$	$::= (\text{init}) \mid (\kappa \text{ arg } \rho \langle V \dots \bullet M \dots \rangle) \mid (\kappa \text{ cond}(M, M, \rho))$		(Cont. code)
$R \in Req$	$::= Req\langle s, sgp, so, \dots \rangle$		(Request)
$so \in SO$	$::= \langle serialised \text{ objects} \rangle$		(Serialised Object)
$sgp \in SO$	$::= \langle serialised \text{ global pointer} \rangle$		(Serialised Object)

Environment Operations:	Store Operations:
$\rho(x) = V$ if $(x V) \in \rho$	$\theta \uplus \{(\alpha V)\} = \theta \cup \{(\alpha V)\}$
$\rho[x \leftarrow V] = (\rho \setminus \{(x V')\}) \cup \{(x V)\}$	with $\alpha \notin DOM(\theta)$
if $(x V') \in \rho$	$\theta(\alpha) = V$ if $(\alpha V) \in \theta$
$\rho[x \leftarrow V] = \rho \cup \{(x V)\}$	$\theta[\alpha := V] = (\theta \setminus \{(\alpha \theta(\alpha))\}) \cup \{(\alpha V)\}$
if $x \notin DOM(\rho)$	
$\rho[x \vec{x} \leftarrow V_1, V_2 \dots] = \rho[x \leftarrow V_1][\vec{x} \leftarrow V_2 \dots]$	$listify(V_1) = (\text{cons } V_1 \text{ nil})$
$\rho[.0 \leftarrow] = \rho$	$listify(V_1, V_2 \dots) = (\text{cons } V_1 listify(V_2 \dots))$
$\rho[x \leftarrow] = \rho[x \leftarrow \text{nil}]$	
$\rho[.x \leftarrow V \dots] = \rho[x \leftarrow listify(V \dots)]$	

---

**Fig. 1.** State Space of the CEKDS-machine

---

---

$\text{Ev}\langle(M\ M_1\dots), \rho, \kappa\rangle \rightarrow_C \text{Ev}\langle M, \rho, (\kappa \text{ arg } \rho \ \langle\bullet, M_1, \dots\rangle)\rangle$	(operator)
$\text{Ev}\langle\lambda\vec{x}.M, \rho, \kappa\rangle \rightarrow_C \text{Ret}\langle\langle\text{cl } \lambda\vec{x}.M, \rho\rangle, \kappa\rangle$	(lambda)
$\text{Ev}\langle c, \rho, \kappa\rangle \rightarrow_C \text{Ret}\langle c, \kappa\rangle$	(constant)
$\text{Ev}\langle x, \rho, \kappa\rangle \rightarrow_C \text{Ret}\langle\rho(x), \kappa\rangle$	(variable)
$\text{Ret}\langle V, (\kappa \text{ arg } \rho \ \langle V_1, \dots, \bullet, M, M_1, \dots\rangle)\rangle \rightarrow_C \text{Ev}\langle M, \rho, (\kappa \text{ arg } \rho \ \langle V_1, \dots, V, \bullet, M_1, \dots\rangle)\rangle$	(operand)
$\text{Ret}\langle V, (\kappa \text{ arg } \rho \ \langle\langle\text{cl } \lambda\vec{x}.M, \rho_1\rangle, V_1, \dots, \bullet\rangle)\rangle \rightarrow_C \text{Ev}\langle M, \rho_1[\vec{x} \leftarrow V_1, \dots, V], \kappa\rangle$	(apply)
$\text{Ret}\langle\langle\text{cl } \lambda\vec{x}.M, \rho\rangle, (\kappa \text{ arg } \rho' \ \langle\bullet\rangle)\rangle \rightarrow_C \text{Ev}\langle M, \rho[\vec{x} \leftarrow ], \kappa\rangle$	(apply0)
$\text{Ev}\langle(\text{if } M\ M_1\ M_2), \rho, \kappa\rangle \rightarrow_C \text{Ev}\langle M, \rho, (\kappa \text{ cond } (M_1, M_2, \rho))\rangle$	(predicate)
$\text{Ret}\langle V, (\kappa \text{ cond } (M, M_1, \rho))\rangle \rightarrow_C \text{Ev}\langle M_1, \rho, \kappa\rangle \quad \text{if } V = \text{false}$	(if else)
$\rightarrow_C \text{Ev}\langle M, \rho, \kappa\rangle \quad \text{if } V \neq \text{false}$	(if then)
$\text{Ret}\langle V, (\kappa \text{ arg } \rho \ \langle\text{callcc}, \bullet\rangle)\rangle \rightarrow_C \text{Ret}\langle\langle\text{co } \kappa\rangle, (\kappa \text{ arg } \rho \ \langle V, \bullet\rangle)\rangle$	(capture)
$\text{Ret}\langle V, (\kappa \text{ arg } \rho \ \langle\langle\text{co } \kappa'\rangle, \bullet\rangle)\rangle \rightarrow_C \text{Ret}\langle V, \kappa'\rangle$	(invoke)
$\text{Ret}\langle V, (\kappa \text{ arg } \rho \ \langle f, V_1, \dots, \bullet\rangle)\rangle \rightarrow_C \text{Ret}\langle\delta(f, V_1, \dots, V), \kappa\rangle$	( $\delta$ )

---

Fig. 2. CEK Machine

The transitions between computational states in Figure 2 deal with the sequential subset of the language that does not involve the store. The transitions extend the CEK-machine transitions [5] by accepting  $n$ -ary applications ( $M_1\ M_2\ \dots$ ) and abstractions ( $\lambda\vec{x}.M$ ) with variable number of arguments. Consequently the continuation code ( $\kappa \text{ arg } \rho \ \langle V_1, \dots, V_n, \bullet, M_{n+2}, \dots \rangle$ ) conveys the following meaning: the first  $n$  components of the application have already been evaluated and their values are  $V_1, \dots, V_n$ ; the next component, i.e., the  $n+1$ th, is being evaluated; the components  $M_{n+2}, \dots$  remain to be evaluated in the environment  $\rho$ .

Figures 3 to 6 deal with site transitions. According to rule (*sequential*), a site can perform a transition if it contains a task that can perform a transition of Figure 2. In rule (*fork*), the primitive fork creates a new task applying the closure received as argument and returns the new task name.

Figure 4 deals with operations related to boxes. The primitive *makeref* allocates a new location  $\alpha$ , stores the value  $V$  in that location, and returns a new box object  $\langle\text{bx } \alpha, s\rangle$ , pointing at the location  $\alpha$  in the current site  $s$ . The primitive *deref* returns the content of the location associated with the box and the primitive *setref* changes the content of its associated location.

Rule (*make gp*) creates global pointers in the same way as (*makeref*) creates boxes. A remote service request has the form  $(\text{rsr } \text{str } \text{gp } M_1 \dots M_n)$ , where *str* must evaluate to a string naming a handler to be called on the site that *gp* is pointing at, with the arguments obtained as values of  $M_1, \dots, M_n$ . The value returned by a remote service request is the distinguished void value, but its effect is to add a *request* in the output queue of the current site. A request is not a first-class object, i.e., it is not part of the set of RValues, but it contains the destination site designated by *gp*, the serialised string, and the serialised values of  $M_1, \dots, M_n$ .

Let us observe that the semantics does not specify the behaviour of the *Serial* and *Deserial* functions. Their purpose is to convert to and from a suitable format for transportation. The composition of these functions returns a result that is an isomorphic copy of its argument. (The specification is available from [15].)

If there is a request in the input queue of a site, rule (*handle*) deserialises its content, and applies the handler associated with *str* to the value indicated by the global pointer and the deserialised values. It is also worth noticing that the incoming request is handled

---


$$\begin{aligned}
& \langle \{ \langle C, \tau \rangle \} \cup T, \theta, s, H, I, O \rangle \\
& \quad \rightarrow_m \langle \{ \langle C_1, \tau \rangle \} \cup T, \theta, s, H, I, O \rangle \quad \text{if } C \rightarrow_C^* C_1 \quad (\text{sequential}) \\
& \langle \{ \langle \text{Ret}(\langle \text{cl } \lambda \vec{x}.M, \rho \rangle, (\kappa \text{ arg } \rho_1 \langle \text{fork}, \bullet \rangle)), \tau \rangle \} \cup T, \theta, s, H, I, O \rangle \\
& \quad \rightarrow_m \langle \{ \langle \text{Ret}(\tau_1, \kappa), \tau \rangle \} \cup \{ \langle \text{Ev}(M, \rho[\vec{x} \leftarrow \cdot], (\text{init})), \tau_1 \rangle \} \cup T, \theta, s, H, I, O \rangle \quad (\text{fork}) \\
& \quad \text{with } \tau_1 \notin \text{FTN}(T) \cup \{\tau\}
\end{aligned}$$

---

**Fig. 3.** Sequential and Parallel Evaluation

---

$$\begin{aligned}
& \langle \{ \langle \text{Ret}(V, (\kappa \text{ arg } \rho \langle \text{makeref}, \bullet \rangle)), \tau \rangle \} \cup T, \theta, s, H, I, O \rangle \\
& \quad \rightarrow_m \langle \{ \langle \text{Ret}(\langle \text{bx } \alpha, s \rangle, \kappa), \tau \rangle \} \cup T, \theta \uplus \{ \langle \alpha V \rangle \}, s, H, I, O \rangle \quad (\text{makeref}) \\
& \langle \{ \langle \text{Ret}(\langle \text{bx } \alpha, s \rangle, (\kappa \text{ arg } \rho \langle \text{deref}, \bullet \rangle)), \tau \rangle \} \cup T, \theta, s, H, I, O \rangle \\
& \quad \rightarrow_m \langle \{ \langle \text{Ret}(\theta(\alpha), \kappa), \tau \rangle \} \cup T, \theta, s, H, I, O \rangle \quad (\text{deref}) \\
& \langle \{ \langle \text{Ret}(V, (\kappa \text{ arg } \rho \langle \text{setref}, \langle \text{bx } \alpha, s \rangle, \bullet \rangle)), \tau \rangle \} \cup T, \theta, s, H, I, O \rangle \\
& \quad \rightarrow_m \langle \{ \langle \text{Ret}(\text{void}, \kappa), \tau \rangle \} \cup T, \theta[\alpha := V], s, H, I, O \rangle \quad (\text{setref})
\end{aligned}$$

---

**Fig. 4.** Boxes

---

$$\begin{aligned}
& \langle \{ \langle \text{Ret}(V, (\kappa \text{ arg } \rho \langle \text{make-gp}, \bullet \rangle)), \tau \rangle \} \cup T, \theta, s, H, I, O \rangle \\
& \quad \rightarrow_m \langle \{ \langle \text{Ret}(\langle \text{gp } \alpha, s \rangle, \kappa), \tau \rangle \} \cup T, \theta \uplus \{ \langle \alpha V \rangle \}, s, H, I, O \rangle \quad (\text{make gp}) \\
& \langle \{ \langle \text{Ret}(V_n, (\kappa \text{ arg } \rho \langle \text{rsr}, \text{str}, \langle \text{gp } \alpha, s \rangle, V_1, \dots, \bullet \rangle)), \tau \rangle \} \cup T, \theta, s_1, H, I, O \rangle \\
& \quad \rightarrow_m \langle \{ \langle \text{Ret}(\text{void}, \kappa), \tau \rangle \} \cup T, \theta, s_1, H, I, \{ \text{Req}(s, \text{sgp}, \text{sstr}, \text{so}_1, \dots, \text{so}_n) \} \cup O \rangle \quad (\text{rsr}) \\
& \quad \text{if } \text{Serial}(s_1, \theta, \langle \langle \text{gp } \alpha, s \rangle, \text{str}, V_1, \dots, V_n \rangle) = \langle \text{sgp}, \text{sstr}, \text{so}_1, \dots, \text{so}_n \rangle \\
& \langle T, \theta, s, H, \{ \text{Req}(s, \text{sgp}, \text{str}, \text{so}_1, \dots, \text{so}_n) \} \cup I, O \rangle \\
& \quad \rightarrow_m \langle \{ \langle \text{Ev}(M, \rho[\vec{x} \leftarrow \theta(\alpha), V_1, \dots, V_n], (\text{init})), \tau \rangle \} \cup T, \theta_1, s, H, I, O \rangle \quad (\text{handle}) \\
& \quad \text{if } H(\text{str}) = \langle \text{cl } \lambda \vec{x}.M, \rho \rangle, \\
& \quad \text{Deserial}(s, \theta, \langle \text{sgp}, \text{so}_1, \dots, \text{so}_n \rangle) = \langle \langle \text{gp } \alpha, s \rangle, V_1, \dots, V_n \rangle, \theta_1 \rangle \\
& \quad \text{with } \tau \notin \text{FTN}(T) \\
& \langle \{ \langle \text{Ret}(V, (\kappa \text{ arg } \rho \langle \text{define-handler}, \text{str}, \bullet \rangle)), \tau \rangle \} \cup T, \theta, s, H, I, O \rangle \\
& \quad \rightarrow_m \langle \{ \langle \text{Ret}(\text{void}, \kappa), \tau \rangle \} \cup T, \theta, s, H \uplus \{ \langle \text{str } V \rangle \}, I, O \rangle \quad (\text{define handler})
\end{aligned}$$

---

**Fig. 5.** Global Pointers and Remote Service Requests

---

$$\begin{aligned}
& \langle \{ \langle \text{Ret}(\text{make-lock}, (\kappa \text{ arg } \rho \langle \bullet \rangle)), \tau \rangle \} \cup T, \theta, s, H, I, O \rangle \\
& \quad \rightarrow_m \langle \{ \langle \text{Ret}(\langle \text{lk } \alpha, s \rangle, \kappa), \tau \rangle \} \cup T, \theta \uplus \{ \langle \alpha \text{ true} \rangle \}, s, H, I, O \rangle \quad (\text{make-lock}) \\
& \langle \{ \langle \text{Ret}(\langle \text{lk } \alpha, s \rangle, (\kappa \text{ arg } \rho \langle \text{lock}, \bullet \rangle)), \tau \rangle \} \cup T, \theta, s, H, I, O \rangle \\
& \quad \rightarrow_m \langle \{ \langle \text{Ret}(\text{void}, \kappa), \tau \rangle \} \cup T, \theta[\alpha := \text{false}], s, H, I, O \rangle \quad \text{if } \theta(\alpha) = \text{true} \quad (\text{acquire}) \\
& \langle \{ \langle \text{Ret}(\langle \text{lk } \alpha, s \rangle, (\kappa \text{ arg } \rho \langle \text{unlock}, \bullet \rangle)), \tau \rangle \} \cup T, \theta, s, H, I, O \rangle \\
& \quad \rightarrow_m \langle \{ \langle \text{Ret}(\text{void}, \kappa), \tau \rangle \} \cup T, \theta[\alpha := \text{true}], s, H, I, O \rangle \quad (\text{release})
\end{aligned}$$

---

**Fig. 6.** Mutual Exclusion

---

---


$$\begin{aligned}
& \{m_1\} \cup \mathcal{W} \\
& \quad \rightarrow_{\mathcal{W}} \{m_2\} \cup \mathcal{W} \quad \text{if } m_1 \rightarrow_m m_2 \quad (\text{site}) \\
& \{ \langle T, \theta, s, H, I, \{ \text{Req}(s, \text{req}) \} \cup O \rangle \cup \mathcal{W} \\
& \quad \rightarrow_{\mathcal{W}} \{ \langle T, \theta, s, H, \{ \text{Req}(s, \text{req}) \} \cup I, O \rangle \cup \mathcal{W} \quad (\text{local request}) \\
& \{ \langle T_0, \theta, s_0, H_0, I_0, \{ \text{Req}(s_1, \text{req}) \} \cup O_0 \rangle, \langle T_1, \theta_1, s_1, H_1, I_1, O_1 \rangle \} \cup \mathcal{W} \\
& \quad \rightarrow_{\mathcal{W}} \{ \langle T_0, \theta, s_0, H_0, I_0, O_0 \rangle, \langle T_1, \theta_1, s_1, H_1, \{ \text{Req}(s_1, \text{req}) \} \cup I_1, O_1 \rangle \} \cup \mathcal{W} \\
& \quad \quad \quad (\text{migrate request}) \\
& \{ \{ \langle \text{Ret}(\langle \text{gp } \alpha, s \rangle, (\kappa \text{ arg } \rho \langle \text{remote}, \text{"host"}, \langle \text{cl } \lambda \vec{x}.M, \emptyset \rangle, \bullet \rangle)), \tau \rangle \} \cup T, \theta, s, H, I, O \rangle \} \cup \mathcal{W} \\
& \quad \rightarrow_{\mathcal{W}} \{ \{ \langle \text{Ret}(\langle \text{gp } \alpha_1, s_1 \rangle, \kappa), \tau \rangle \} \cup T, \theta, s, H, I, O, m \} \cup \mathcal{W} \quad (\text{remote}) \\
& \quad \text{with } m = \{ \langle \text{Ev}(\langle M, \emptyset[\vec{x} \leftarrow \langle \text{gp } \alpha, s \rangle], (\text{init})), \tau \rangle, \theta_1, s_1, \emptyset, \emptyset, \emptyset \rangle \\
& \quad \text{with } \theta_1 = \{ (\alpha_1 \text{ nil}) \}, s_1 \notin \text{FSN}(\mathcal{W}) \cup \{s\}
\end{aligned}$$


---

Fig. 7. Other Distribution Aspects

*asynchronously* and that the handler is executed in a new thread. Optimistic active-messages [25] show that the thread could be created lazily [13] and still produce good performance. According to rule (*define handler*), the primitive *define-handler* allows a program to add new handlers in a given site.

Idealised NeXeme offers some primitives to implement critical sections. A *lock* is a datastructure that can be acquired by at most one task using the primitive *lock*, and that can be released by the primitive *unlock*. In Figure 6, a busy-wait implementation of locks relies on the fact that transition rules are atomic.

Figure 7 displays transitions between worlds. Rule (*site*) shows that a world can perform a transition, if there is a site which can perform a transition of Figures 3 to 6. In addition, rules (*local request*) and (*migrate request*) deal with request transfers, and rule (*remote*) is used to create a new site in a world.

Creation of a remote site is always a delicate problem in a distributed system. Frequently, libraries for distribution assume an SPMD approach and use operating-system primitives to startup remote processes. We decided not to model NeXeme very closely, as the semantics would have been unsatisfactorily complicated. Instead, Idealised NeXeme provides a primitive to create a new remote site and to initiate a computation on that site. This initial computation is essential, because it allows the user to install initial handlers for remote service requests. In addition, it makes it possible to synchronise both sites so that the initial site is kept informed when the new site is initialised and is ready to perform some processing.

A remote site is created by the primitive *remote* which takes the following arguments: the name of a host, a closure representing the initial computation to be done, and a global pointer that can be used by the remote site to notify the end of initialisation by a remote service request. The value returned by the primitive *remote* is a global pointer pointing at the new remote site. In order to be implementable in languages without closures, like C, rule (*remote*) requires the closure to have an empty environment. As a result, the closure can be encoded as a pointer to its entry point. Let us observe that the remote site should be executed on “host”.

Having defined a transition relation between worlds, we can now define an evaluation relation associating programs with values. Detecting termination in a distributed system is a difficult matter in the general case. In Idealised NeXeme, we assume there exists a primitive *shutdown* that the program has to call on the origin site; it terminates the computation and returns its argument.

**Definition 1 (Eval)** Let  $M$  be a closed term. An *initial task* is of the form  $t_i =$

$\langle \text{Ev}(M, \emptyset, ((\text{init}))), \tau \rangle$ , while a *final task* is of the form  $t_f = \langle \text{Ret}(V, (\kappa \text{ arg } \rho \langle \text{shutdown}, \bullet \rangle)), \tau' \rangle$ . We have that:

$$\text{eval}(M) = V \text{ if } \langle \{ t_i \}, \emptyset, s_0, \emptyset, \emptyset, \emptyset \rangle \rightarrow_{\mathcal{W}} \{ \langle \{ t_f \} \cup T, \theta, s_0, H, I, O \rangle \} \cup \mathcal{W}.$$

□

## 4 The Implementation

### 4.1 General Description

Figure 8 displays the organisation of the NeXeme implementation. As Nexus is multi-threaded, we had to adopt a thread-safe garbage collector. The public domain Boehm-Weiser's [11] garbage collector supports various OS-level threads and is also part of the PPCR portable common runtime system [26]. On each platform, Nexus is recompiled against the garbage collector (abbreviated gc) and the suitable thread package. The executable is generated by linking the Scheme system with the resulting libraries Nexus, gc, and threads.

---

Utilities
NeXeme
Functional Nexus
Foreign Interface
Nexus Threads Scheme
GC

---

**Fig. 8.** Organisation of NeXeme

Primitives of the Nexus library are made accessible through a foreign interface definition. Let us note that some Nexus primitives require procedures as argument (for instance, thread creation or callbacks for handlers). In order to integrate properly Nexus with Scheme, we modified Nexus to support callbacks to Scheme functions. While the foreign interface defines a similar API as the Nexus library, the “functional Nexus” layer provides a more functional interface to Nexus; for instance, results are returned by functions, errors are signalled by exceptions, and memory is managed automatically. The NeXeme layer offers a functional version of remote service requests as described in Section 3. Finally, a library defining a set of utilities provides other paradigms for distribution like futures, communication channels, or farm processing.

### 4.2 Distributed Garbage Collection

NeXeme has a distributed garbage collector that takes care of memory management automatically. Our working hypotheses, provided by Nexus, are a reliable message-passing and a FIFO ordering<sup>4</sup> of messages between two sites.

Each site relies on a thread safe, conservative, mark and sweep garbage collector [11]; conservativeness is required as Scheme data are passed to Nexus, written in C, and are pointed by Nexus data structures. In addition, NeXeme maintains two tables for each site. The *exit table*<sup>5</sup> associates each global pointer with the number of distinct remote copies of this pointer originating from the site. The *entry table* contains all global pointers received by a site, except those that point at itself. The exit table, but not the entry table, is a root of the local garbage collector. The role of the distributed collector is to update counters in a safe and consistent way. To this end, it relies on two types of *control* messages, called “*decrement(gp)*” and “*increment-decrement(gp, s)*”, as described below.

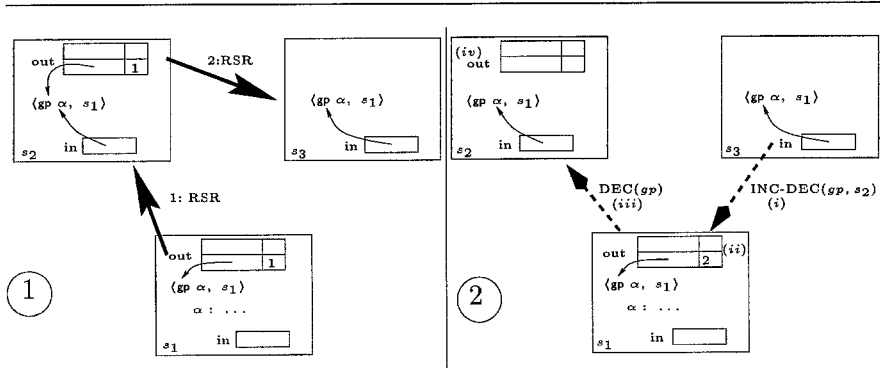
<sup>4</sup> Nexus, as opposed to Idealised NeXeme, allows requests to be handled in a non-threaded manner; in such a case, Nexus guarantees the FIFO ordering of requests between two sites.

<sup>5</sup> Different names are found in the literature, such as entry/exit items, scions/stubs [17].

Reference counters are updated according to the *diffusion tree* [16] of global pointers. The first time a global pointer  $gp$  is serialised, an entry is added into the exit table of the current site with a counter set to 1; afterwards, for every serialisation, this counter is incremented. Symmetrically, the first time a global pointer  $gp$  is deserialised, it is added to the entry table (if it points at a remote host); if it is already present, NeXeme sends a decrement message for this global pointer, “ $decrement(gp)$ ” to the site that sent the remote service request. When a site receives a message “ $decrement(gp)$ ”, the counter of  $gp$  in the exit table is decremented. Once the counter reaches zero, the global pointer may be removed from the table; only then, the pointer itself and the local data may be reclaimed, if no longer accessible.

In Nexus, a global pointer contains the site it is pointing at, and not the site it is arriving from. Therefore, once a global pointer  $gp$  pointing at a site  $s_1$  becomes inaccessible on a site  $s_2$ ,  $s_2$  sends a message “ $decrement(gp)$ ” to  $s_1$ . This naive implementation of the reference counter technique is sound if causality is preserved in the system [16]. This is ensured by a reorganisation of the diffusion tree, as explained in the following scenario illustrated in Figure 9.

Let us consider that  $gp$ , a global pointer pointing to address  $\alpha$  in  $s_1$ , is migrated to  $s_2$ , and then migrated to  $s_3$ . Figure 9 part 1 shows that the reference counters in  $s_1$  and  $s_2$  are equal to 1, meaning that there are  $1 + 1 = 2$  active remote references of  $gp$ . Once  $s_3$  deserialises  $gp$ , the global pointer received from  $s_2$ , a reorganisation can be initiated, as illustrated in Figure 9 part 2. (i) Site  $s_3$  sends an “ $increment-decrement(gp, s_2)$ ” message to  $s_1$ ; (ii) when the message is received by  $s_1$ , the counter for  $gp$  is incremented on  $s_1$ ; (iii) afterwards, a “ $decrement(gp)$ ” message is sent to  $s_2$ ; (iv) when the message is received, the counter for  $gp$  is decremented on  $s_2$ .



**Fig. 9.** Distributed Reference Counters (1) after Remote Service Requests  
(2) after Control Messages

Race conditions are avoided between  $s_1$  and  $s_3$  by giving priority to “increment-decrement” over “decrement” messages. As a result, we have the following invariant. For every global pointer  $gp$  pointing to site  $s$ :

$$\sum_{site} counter(gp) = \sum_{site \neq s} in(gp) + \sum Value(message \text{ for } gp)$$

with  $Value(decrement(gp)) = 1$  and  $Value(increment-decrement(gp, s_i)) = 0$ .

The simplicity and portability of the solution is unfortunately counter-balanced by its inability to collect distributed cycles. Our approach differs from “Indirect Reference Counting” [16] because it can reclaim “zombie” pointers, i.e.,  $gp$  can be freed on  $s_2$  even though it remains active on  $s_3$ ; other techniques to reduce the length of indirection chains



are described in [17]. Several optimisations are possible. First, control messages of the same type may be grouped in a single message. Second, an “*increment-decrement(gp, s<sub>2</sub>)*” to be sent to  $s_1$ , followed by a “*decrement(gp)*” to the same destination, may be replaced by a “*decrement(gp)*” to  $s_2$ .

#### 4.3 Utilities

The utilities layer offers a number of services that we briefly describe in this section. *Remote invocation of a procedure* sends a locally-available closure to a remote site where it is invoked on some arguments; as in a remote service request, no result is returned by this operation. Not only does *remote invocation of a function* transfer and call a function on some arguments, but also it returns the result to the site that initiated the operation. Communication channels [21, 9], and futures [12] are also provided.

“Farm-processing” versions of the various remote call utilities are provided that can be used when a list of global pointers is passed as an argument; these versions map the corresponding operation onto each element of the list. NeXeme is also able to detect when a site becomes idle so that lazy task creation [13] and task stealing can be implemented easily.

### 5 Discussion and Related Work

The Nexus [7] philosophy is derived from *active messages* [24], where each message contains at its head the address of a user-level handler executed on message arrival. Nexus has also been used to build several other distributed languages. In Fortran M [6], program modules, called processes, communicate via explicitly-declared channels. The language nPerl [8] is a distributed extension of Perl which provides a remote procedure call; unlike NeXeme, higher-order functions and continuations are not available. Compositional C++ [4] supports global pointers that can refer to objects on remote machines, and defines usual operations on them; also, CC+ offers single assignment objects. Currently, NeXeme uses Nexus 3.0; Nexus 4.0 provides EZ Nexus a more abstract interface to remote service requests, which is similar to the one described here.

Boxes allow the programmer to define mutable data structures, for example cyclic data structures. If a cyclic data structure is passed as an argument to a remote service request, an isomorphic copy is generated on the remote site. NeXeme does not prevent the user from mutating each copy separately, as they are now independent entities. However, using global pointers, any form of consistency can be programmed on top of NeXeme, such as simulated shared memory [14] or causally-coherent memory [18].

There exist other distributed implementations of functional languages with *explicit* notations for parallelism or distribution. DMeroon [18] offers a distributed memory model which enforces coherence based on causality, while ICSLAS [19] adds to DMeroon transparent remote execution. Kali Scheme [3] also offers a form of active messages, on top of which higher-level primitives are defined; the semantics of Kali Scheme could easily be defined in terms of Idealised NeXeme. Distributed versions of SML-style languages are usually based on communications channels [21, 9].

### 6 Conclusion

This paper presents NeXeme, a distributed dialect of Scheme, based on remote service requests and global pointers provided by the library for distribution Nexus. This paper describes the semantics but also the implementation of this language. The functional interface to remote service requests of NeXeme is a perfect abstraction to build a Scheme with futures [14]; in addition, NeXeme is used for programming multimedia distributed applications.

### 7 Acknowledgement

Hans Boehm, Carl Hauser, Carl Kesselman, Andreas Kind, Christian Queinnec, Manuel Serrano, and Steve Tuecke gave useful comments and advice for the implementation of NeXeme.

## Bibliography

- [1] A. D. Birell and B. J. Nelson. Implementing Remote Procedure Calls. *ACM Transactions on Computer Systems*, 2(1):39–59, February 1984.
- [2] L. Cardelli. A Language with Distributed Scope. In *Proceedings of the Twenty Second Annual ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages*, pages 286–297, San-Francisco, California, January 1995.
- [3] H. Cejtin, S. Jagannathan, and R. Kelsey. Higher-Order Distributed Objects. *ACM Transactions on Programming Languages and Systems*, 17(5):704–739, 1995.
- [4] K. M. Chandy and C. Kesselman. CC++: A declarative concurrent object oriented programming notation. In *Research Directions in Object Oriented Programming*, pages 281–313. MIT Press, 1993.
- [5] M. Felleisen and D. P. Friedman. Control Operators, the SECD-Machine and the  $\lambda$ -Calculus. In M. Wirsing, editor, *Formal Description of Programming Concepts III*, pages 193–217, Amsterdam, 1986. Elsevier Science Publishers B.V. (North-Holland).
- [6] I. Foster and K.M. Chandy. Fortran M: A Language for Modular Parallel Programming. *J. of Parallel and Distributed Computing*, 26(1):24–35, 1995.
- [7] I. Foster, C. Kesselman, and S. Tuecke. The Nexus Approach to Integrating Multithreading and Communication. *J. of Parallel and Distributed Computing*, 37:70–82, 1996.
- [8] I. Foster and R. Olson. A Guide to Parallel and Distributed Programming in nPerl. Mathematics and Computer Science Division, October 1995.
- [9] A. Giacalone, P. Mishra, and S. Prasad. Facile: A Symmetric Integration of Concurrent and Functional Programming. *Int. J. of Parallel Programming*, 18(2):121–160, 1989.
- [10] J. Gosling, G. Steele, and B. Joy. *The Java Language Specification*. Addison-Wesley, 1996.
- [11] H.-J. Boehm and M. Weiser. Garbage Collection in an Uncooperative Environment. *Software – Practice and Experience*, 18(9):807–820, 1988.
- [12] R. Halstead, Jr. Parallel Symbolic Computing. *IEEE Computer*, pages 35–43, Aug. 1986.
- [13] E. Mohr, D. A. Kranz, and R. H. Halstead. Lazy Task Creation : a Technique for Increasing the Granularity of Parallel Programs. In *Proceedings of the 1990 ACM Conference on Lisp and Functional Programming*, pages 185–197, June 1990.
- [14] L. Moreau. Correctness of a Distributed-Memory Model for Scheme. In *Second International Europar Conference (EURO-PAR'96)*, LNCS 1123, pages 615–624, 1996.
- [15] NeXeme Home Page. <http://www.ecs.soton.ac.uk/~lavm/NeXeme>.
- [16] J. M. Piquer. Indirect Distributed Garbage Collection: Handling Object Migration. *ACM Transactions on Programming Languages and Systems*, 18(5):615–647, 1996.
- [17] D. Plainfossé and M. Shapiro. A survey of distributed garbage collection techniques. In *Int. Workshop on Memory Management (IWMM95)*, LNCS 986, pages 211–249, 1995.
- [18] C. Queinnec. DMEROON: a Distributed Class-based Causally-coherent Data Model: Preliminary Report. In *Parallel Symbolic Languages and Systems.*, LNCS 1068, 1995.
- [19] C. Queinnec and D. De Roure. Design of a Concurrent and Distributed Language. In *Parallel Symbolic Computing: Languages, Systems and Applications*, LNCS 748, pages 234–259, Boston, Massachusetts, October 1992.
- [20] J. Rees and W. Clinger, editors. Revised<sup>4</sup> Report on the Algorithmic Language Scheme. *Lisp Pointers*, 4(3):1–55, July–September 1991.
- [21] J. H. Reppy. *Higher-Order Concurrency*. PhD thesis, Department of Computer Science, Cornell University, Ithaca, New York, June 1992.
- [22] Sun Microsystems. Java Remote Method Invocation Specification, November 1996.
- [23] P. van Roy, S. Haridi, P. Brand, and G. Smolka et al. Mobile Objects in Distributed Oz. Technical report, Swedish Institute of Computer Science, December 1996.
- [24] T. von Eicken, D. E. Culler, S. Goldstein, and K. Schausser. Active Messages: a Mechanism for Integrated Communication & Computation. In *Proceedings of the 19th symposium on Computer Architecture*, pages 256–266, 1992.
- [25] D. A. Wallach, W. C. Hsieh, K. Johnson, M. F. Kaashoek, et al. Optimistic active messages: A mechanism for scheduling communication with computation. In *5th ACM SIGPLAN Symp. on Principles and Practice of Parallel Programming (PPoPP '95)*, 1995.
- [26] M. Weiser, A. Demers, and C. Hauser. The Portable Common Runtime Approach to Interoperability. In *ACM Symp. on Operating System Principles*, pages 114–122, Dec. 1989.