

# *Putting logic-based distributed systems on stable grounds*

TOM J. AMELOOT\* and JAN VAN DEN BUSSCHE

*Computer Science, Hasselt, Limburg, Belgium*

(e-mail: tom.ameloot@uhasselt.be, jan.vandenbussche@uhasselt.be)

WILLIAM R. MARCZAK, PETER ALVARO and JOSEPH M. HELLERSTEIN

*Computer Science, Berkeley, California, USA*

(e-mail: wrm@cs.berkeley.edu, palvaro@cs.berkeley.edu, hellerstein@cs.berkeley.edu)

*submitted 5 September 2012; revised 12 November 2013, 24 April 2015; accepted 16 July 2015*

---

## Abstract

In the Declarative Networking paradigm, Datalog-like languages are used to express distributed computations. Whereas recently formal operational semantics for these languages have been developed, a corresponding declarative semantics has been lacking so far. The challenge is to capture precisely the amount of nondeterminism that is inherent to distributed computations due to concurrency, networking delays, and asynchronous communication. This paper shows how a declarative, model-based semantics can be obtained by simply using the well-known stable model semantics for Datalog with negation. We show that the model-based semantics matches previously proposed formal operational semantics.

**KEYWORDS:** Dedalus, Datalog, stable model semantics, distributed system, asynchronous communication

---

## 1 Introduction

Cloud environments have emerged as a modern way to store and manipulate data (Zhang *et al.* 2010; Cavage 2013). For our purposes, a cloud is a distributed system that should produce output as the result of some computation. We use the common term “node” as a synonym for an individual computer or server in a network.

In recent years, logic programming has been proposed as an attractive foundation for distributed and cloud programming, building on work in declarative networking (Loo *et al.* 2009). The essential idea in declarative networking, is that the programmer uses a high-level declarative language (like Datalog) to specify only what has to happen, and not exactly how. For example, the programmer could specify only that certain messages are generated in reply to other messages; the exact technical details

\* T.J. Ameloot is a Postdoctoral Fellow of the Research Foundation – Flanders (FWO).

to send (and possibly resend) messages over transmission protocols are filled in by some runtime engine. This frees the programmer from thinking in low-level terms that distract from the actual meaning of the specific program at hand. In particular, complex distributed algorithms and protocols can be expressed in relatively few lines of code (Jim 2001; Alvaro *et al.* 2009; Hellerstein 2010b). Besides the interest in declarative networking, we are also seeing a more general resurgence of Datalog (with negation) (de Moor *et al.* 2011; Huang *et al.* 2011). Moreover, issues related to data-oriented distributed computing are receiving attention at database theory conferences (Hellerstein 2010a; Abiteboul *et al.* 2011; Ameloot *et al.* 2011; Ameloot and Van den Bussche 2012; Zinn *et al.* 2012).

One of the latest languages proposed in declarative networking is Dedalus (Alvaro *et al.* 2009; Hellerstein 2010b; Alvaro *et al.* 2011), a Datalog-inspired language that has influenced other recent language designs for distributed and cloud computing such as Webdamlog (Abiteboul *et al.* 2011) and Bloom (Alvaro *et al.* 2011).

### 1.1 Model-based semantics

In this paper, we describe the meaning of distributed Datalog programs using a model-based semantics. This approach contrasts with most previous work in declarative networking, where the meaning of programs was typically described with an operational semantics (Deutsch *et al.* 2006; Navarro and Rybalchenko 2009; Grumbach and Wang 2010; Ameloot *et al.* 2011), with a few exceptions (Lobo *et al.* 2012; Ma *et al.* 2013).

There are several important motivations for a model-based semantics of a distributed program. First, we can better separate the program structure, i.e., the rules, from the (distributed) implementation that may change over time. For example, consider rules that generate messages. These rules *can* be implemented with asynchronous communication, but how we evaluate them across machines is eventually just a physical performance decision. Said differently, the point of message rules is not to model a physical phenomenon, but rather to admit a wider array of physical implementations than a local evaluation strategy. Model-based interpretations of a program admit all such implementations, and can perhaps suggest some new ones. Second, we can investigate the *need* for time: we can think about when temporal delay is needed for expressivity, rather than when it is imposed upon us by some implementation detail like physical separation of nodes. In this context, we mention the CRON conjecture by Hellerstein, that relates causality on messages to the nature of the computations in which those messages participate (Hellerstein 2010b; Ameloot and den Bussche 2014). We elaborate on causality below.

Concretely, our approach will be to model a distributed program with Datalog under the stable model semantics (Gelfond and Lifschitz 1988) because this semantics is widely used in logic programming. Following the language Dedalus (Alvaro *et al.* 2009; Hellerstein 2010b; Alvaro *et al.* 2011), we express the functionality of the distributed program with three kinds of rules: “deductive rules” for local computation, “inductive rules” for persisting memory across local computation steps,

and, “asynchronous rules” for representing message sending. The asynchronous rules will nondeterministically choose the arrival times of messages (Krishnamurthy and Naqvi 1988; Saccà and Zaniolo 1990).

However, using only the above rules is not sufficient, as this still allows stable models that express undesirable computations, where messages can be sent “into the past”. Therefore, each program is augmented with a set of rules that express *causality* on the messages. Causality stands for the physical constraint that an effect can only happen after its cause. Applied to message delivery, this intuitively means that a sent message can only be delivered in the future, not in the past. The rules for causality reason from the perspective of the local times of each node, which is a justified approach since there is no common “global clock” in a distributed environment (Attiya and Welch 2004). As a second improvement, we also introduce rules to ensure that only a finite number of messages arrive at each local step of a node, as occurs in a real distributed system. Applying the stable model semantics to the augmented Datalog programs constitutes our modeling of a distributed (Datalog) program.

On another note, it is already well known that for finite input domains, the combination of Datalog and stable model semantics allows for expressing all problems in NP (Marek and Truszczynski 1999). However, it is not yet clear what can be represented when infinite input domains are considered. From this perspective, our work demonstrates that the stable model semantics is indeed also suitable for modeling distributed programs, whose execution is unbounded in time. Here, time would be provided as an infinite input.

## 1.2 Correctness

As we have motivated above, our goal is to describe the workings of a distributed system declaratively, so that new insights can emerge from this perspective. Hence, it is important to verify that the model-based semantics really corresponds to the execution of a distributed program.

To this end, we additionally formalize the execution of a distributed Datalog program by means of an *operational* semantics (Deutsch et al. 2006; Navarro and Rybalchenko 2009; Grumbach and Wang 2010; Ameloot et al. 2011). This second semantics is defined as a transition system. The transition system is infinite because nodes run indefinitely and keep sending messages. In addition, the transition system is highly nondeterministic, because nodes work concurrently and messages can be delayed.

We establish rigorously a correspondence between the features of the operational semantics and the features of the proposed model-based semantics. To formulate our result, we describe each operational execution by a structure that we call a *trace*, which includes for each node in the network the detailed information about the local steps it has performed and about the messages it has sent and received. For our distributed Datalog programs, we show that such operational traces correspond to the set of stable models.

### 1.3 Outline

This paper is organized as follows. First, Section 2 discusses related work. Section 3 gives preliminaries. Next, Section 4 represents distributed Datalog programs under the model-based semantics; this section is based on Dedalus, a Datalog-like language. Section 5 justifies the intuitions of the model-based semantics by establishing an equivalence with an operational semantics. Section 6 finishes with the conclusion.

## 2 Related work

The work of Lobo *et al.* (2012) is closely related to our work. For a Dedalus-inspired language, they give a model-theoretic semantics based on answer set programming, i.e., stable models. To define this semantics, they syntactically translate the rules of their language to Datalog, where all literals are given an explicit location and time variable, to represent the data that each node has during each local time. This translation resembles the model-theoretic semantics for distributed Datalog programs in this paper. To enforce natural execution properties in their semantics, like causality, Lobo *et al.* specify auxiliary rules in the syntactical translation. The work of Lobo *et al.* (2012) does not yet mention the connection between the model-theoretic semantics and desired executions of a distributed system, i.e., an operational semantics.

Extending the work of Lobo *et al.*, the work of Ma *et al.* (2013) formalizes a distributed system as a composition of I/O automata (Lynch 1996). An operational execution of such a system is a sequence of valid transitions, called a trace. Global properties of the system can be analyzed by translating it into a logic program, to which an answer set solver can be applied. Ma *et al.* mention that operational traces of the system correspond to answer sets of the logic program, and that this provides a formal foundation for the analysis tools based on answer set programming. Thus, the work of Ma *et al.* (2013) indicates a practical benefit of having a correspondence between a declarative and operational semantics for languages used in declarative networking. As mentioned above, we also establish a similar correspondence in the current paper, for our distributed Datalog programs. We note, however, a few differences between our work and that of Ma *et al.* First, in the work of Ma *et al.*, the message buffer of a node has a maximum size. In our operational semantics, the buffers are unbounded. Moreover, Ma *et al.* construct their logic programs for a fixed range of timestamps. In our declarative, model-based semantics, time is given as an infinite input to a Datalog program whose rules are independent of a fixed time range. Lastly, our work devotes much attention to rigorously showing the correspondence between the declarative and operational semantics, whereas this is not elaborated in the work of Ma *et al.*

Also in the setting of distributed systems, Interlandi *et al.* (2013) give a Dedalus-inspired language for describing *synchronous* systems. In such systems, the nodes of the network proceed in rounds and the messages cannot be arbitrarily delayed. During each round, the nodes share the same global clock. Interlandi *et al.* specify an operational semantics for their language, based on relational transducer networks

(Ameloot *et al.* 2013). They also show that this operational semantics coincides with a model-theoretic semantics of a single holistic Datalog program. It should be noted that Lobo *et al.* (2012), and the current paper, deal with *asynchronous* systems, that in general pose a bigger challenge for a distributed program to be correct, i.e., the program should remain unaffected by nondeterministic effects caused by message delays.

An area of artificial intelligence that is closely related to declarative networking is that of programming multi-agent systems in declarative languages. The knowledge of an agent can be expressed by a logic program, which also allows for non-monotone reasoning, and agents update their knowledge by modifying the rules in these logic programs (Leite *et al.* 2002; Nigam and Leite 2006; Leite and Soares 2007). The language LUPS (Alferes *et al.* 2002) was designed to specify such dynamic updates to logic programs, and LUPS is also a declarative language itself. After applying a sequence of updates specified in LUPS, the semantics of the resulting logic program can be defined in an inductive way. But an interesting connection to this current work, is that the semantics can also be given by first syntactically translating the original program and its updates into a single normal logic program, after which the stable model semantics is applied (Alferes *et al.* 2002). It should be noted however that in this second semantics, there is no modeling of causality or the sending of messages.

Of course, logic programming is not the only means for specifying a (distributed) system. For example, in the area of formal methods, logic-based languages like TLA (Lamport 2000a), Z (Woodcock and Davies 1996), and Event-B (Abrial 2010) can be used to specify various distributed algorithms. Specifications written in these languages can also be automatically checked for correctness.

Although we work within the established setting of declarative networking (Loo *et al.* 2009), the scientific debate on the merits of Datalog versus other formalisms for programming distributed systems remains open. It seems desirable to have an analysis of how features of Datalog relate to the features of other languages for formal specification, e.g. (Woodcock and Davies 1996; Lamport 2000a; Abrial 2010), both on the syntactical and the semantical level. However, a deep understanding of the other languages would be needed. Moreover, one may expect that features of Datalog will in general not map naturally to features of the other languages. Hence, we consider such a comparison to be a separate research project, outside the scope of the current paper.

### 3 Preliminaries

#### 3.1 Database basics

A *database schema*  $\mathcal{D}$  is a finite set of pairs  $(R, k)$  where  $R$  is a *relation name* and  $k \in \mathbb{N}$  its associated *arity*. A relation name occurs at most once in a database schema. We often write  $(R, k)$  as  $R/k$ .

We assume some infinite universe **dom** of atomic data values. A *fact*  $\mathbf{f}$  is a pair  $(R, \bar{a})$ , often denoted as  $R(\bar{a})$ , where  $R$  is a relation name and  $\bar{a}$  is a tuple of values

over **dom**. For a fact  $R(\bar{a})$ , we call  $R$  the *predicate*. We say that a fact  $R(a_1, \dots, a_k)$  is *over* database schema  $\mathcal{D}$  if  $R/k \in \mathcal{D}$ . A database *instance*  $I$  over  $\mathcal{D}$  is a set of facts over  $\mathcal{D}$ . For a subset  $\mathcal{D}' \subseteq \mathcal{D}$ , we write  $I|_{\mathcal{D}'}$  to denote the subset of facts in  $I$  whose predicate is a relation name in  $\mathcal{D}'$ . We write  $adom(I)$  to denote the set of values occurring in facts of  $I$ .

### 3.2 Datalog with negation

We recall Datalog with negation (Abiteboul *et al.* 1995), abbreviated Datalog<sup>−</sup>. We assume the standard database perspective, where a Datalog<sup>−</sup> program is evaluated over a given set of facts, i.e., where these facts are not part of the program itself.

Let **var** be a universe of *variables*, disjoint from **dom**. An *atom* is of the form  $R(u_1, \dots, u_k)$  where  $R$  is a relation name and  $u_i \in \mathbf{var} \cup \mathbf{dom}$  for each  $i = 1, \dots, k$ . We call  $R$  the *predicate*. If an atom contains no data values, we call it *constant-free*. A *literal* is an atom or an atom with “ $\neg$ ” prepended. A literal that is an atom is called *positive* and otherwise it is called *negative*.

It will be technically convenient to use a slightly unconventional definition of rules. Formally, a Datalog<sup>−</sup> rule  $\varphi$  is a triple

$$(head_\varphi, pos_\varphi, neg_\varphi),$$

where  $head_\varphi$  is an atom;  $pos_\varphi$  and  $neg_\varphi$  are sets of atoms; and, the variables in  $\varphi$  all occur in  $pos_\varphi$ . This last condition is called *safety*. The components  $head_\varphi$ ,  $pos_\varphi$  and  $neg_\varphi$  are called respectively the *head*, the *positive body atoms* and the *negative body atoms*. We refer to  $pos_\varphi \cup neg_\varphi$  as the *body atoms*. Note,  $neg_\varphi$  contains just atoms, not negative literals. Every Datalog<sup>−</sup> rule  $\varphi$  must have a head, whereas  $pos_\varphi$  and  $neg_\varphi$  may be empty. If  $neg_\varphi = \emptyset$  then  $\varphi$  is called *positive*.

A rule  $\varphi$  may be written in the conventional syntax. For instance, if  $head_\varphi = T(u, v)$ ,  $pos_\varphi = \{R(u, v)\}$  and  $neg_\varphi = \{S(v)\}$ , with  $u, v \in \mathbf{var}$ , then we can write  $\varphi$  as

$$T(u, v) \leftarrow R(u, v), \neg S(v).$$

The specific ordering of literals to the right of the arrow has no significance in this paper.

The set of variables of  $\varphi$  is denoted  $vars(\varphi)$ . If  $vars(\varphi) = \emptyset$  then  $\varphi$  is called *ground*, in which case  $\{head_\varphi\} \cup pos_\varphi \cup neg_\varphi$  is a set of facts.

Let  $\mathcal{D}$  be a database schema. A rule  $\varphi$  is said to be *over schema*  $\mathcal{D}$  if for each atom  $R(u_1, \dots, u_k) \in \{head_\varphi\} \cup pos_\varphi \cup neg_\varphi$ , we have  $R/k \in \mathcal{D}$ . A Datalog<sup>−</sup> *program*  $P$  over  $\mathcal{D}$  is a set of (safe) Datalog<sup>−</sup> rules over  $\mathcal{D}$ . We write  $sch(P)$  to denote the smallest database schema that  $P$  is over; note,  $sch(P)$  is uniquely defined. We define  $idb(P) \subseteq sch(P)$  to be the database schema consisting of all relations in rule-heads of  $P$ . We abbreviate  $edb(P) = sch(P) \setminus idb(P)$ .<sup>1</sup>

Any database instance  $I$  over  $sch(P)$  can be given as *input* to  $P$ . Note,  $I$  may already contain facts over  $idb(P)$ .<sup>2</sup> Let  $\varphi \in P$ . A *valuation for*  $\varphi$  is a total function

<sup>1</sup> The abbreviation “idb” stands for “intensional database schema” and “edb” stands for “extensional database schema” (Abiteboul *et al.* 1995).

<sup>2</sup> The need for this will become clear in Section 5.

$V : \text{vars}(\varphi) \rightarrow \mathbf{dom}$ . The application of  $V$  to an atom  $R(u_1, \dots, u_k)$  of  $\varphi$ , denoted  $V(R(u_1, \dots, u_k))$ , results in the fact  $R(a_1, \dots, a_k)$  where for each  $i \in \{1, \dots, k\}$  we have  $a_i = V(u_i)$  if  $u_i \in \mathbf{var}$  and  $a_i = u_i$  otherwise. In words: applying  $V$  replaces the variables by data values and leaves the old data values unchanged. This is naturally extended to a set of atoms, which results in a set of facts. Valuation  $V$  is said to be satisfying for  $\varphi$  on  $I$  if  $V(\text{pos}_\varphi) \subseteq I$  and  $V(\text{neg}_\varphi) \cap I = \emptyset$ . If so,  $\varphi$  is said to derive the fact  $V(\text{head}_\varphi)$ .

### 3.2.1 Positive and semi-positive

Let  $P$  be a Datalog<sup>-</sup> program. We say that  $P$  is *positive* if all rules of  $P$  are positive. We say that  $P$  is *semi-positive* if for each rule  $\varphi \in P$ , the atoms of  $\text{neg}_\varphi$  are over  $\text{edb}(P)$ . Note, positive programs are semi-positive.

We now give the semantics of a semi-positive Datalog<sup>-</sup> program  $P$  (Abiteboul et al. 1995). First, let  $T_P$  be the *immediate consequence operator* that maps each instance  $J$  over  $\text{sch}(P)$  to the instance  $J' = J \cup A$  where  $A$  is the set of facts derived by all possible satisfying valuations for the rules of  $P$  on  $J$ .

Let  $I$  be an instance over  $\text{sch}(P)$ . Consider the infinite sequence  $I_0, I_1, I_2$ , etc, inductively defined as follows:  $I_0 = I$  and  $I_i = T_P(I_{i-1})$  for each  $i \geq 1$ . The *output of  $P$  on input  $I$* , denoted  $P(I)$ , is defined as  $\bigcup_j I_j$ ; this is the *minimal fixpoint* of the  $T_P$  operator. Note,  $I \subseteq P(I)$ . When  $I$  is finite, the fixpoint is finite and can be computed in polynomial time according to data complexity (Vardi 1982).

### 3.2.2 Stratified semantics

We now recall the stratified semantics for a Datalog<sup>-</sup> program  $P$  (Abiteboul et al. 1995). As a slight abuse of notation, here we will treat  $\text{idb}(P)$  as a set of only relation names (without associated arities). First,  $P$  is called *syntactically stratifiable* if there is a function  $\sigma : \text{idb}(P) \rightarrow \{1, \dots, |\text{idb}(P)|\}$  such that for each rule  $\varphi \in P$ , having some head predicate  $T$ , the following conditions are satisfied:

- $\sigma(R) \leq \sigma(T)$  for each  $R(\bar{u}) \in \text{pos}_\varphi|_{\text{idb}(P)}$ ;
- $\sigma(R) < \sigma(T)$  for each  $R(\bar{u}) \in \text{neg}_\varphi|_{\text{idb}(P)}$ .

For  $R \in \text{idb}(P)$ , we call  $\sigma(R)$  the *stratum number* of  $R$ . For technical convenience, we may assume that if there is an  $R \in \text{idb}(P)$  with  $\sigma(R) > 1$  then there is an  $S \in \text{idb}(P)$  with  $\sigma(S) = \sigma(R) - 1$ . Intuitively, function  $\sigma$  partitions  $P$  into a sequence of semi-positive Datalog<sup>-</sup> programs  $P_1, \dots, P_k$  with  $k \leq |\text{idb}(P)|$  such that for each  $i = 1, \dots, k$ , the program  $P_i$  contains the rules of  $P$  whose head predicate has stratum number  $i$ . This sequence is called a *syntactic stratification* of  $P$ . We can now apply the *stratified semantics* to  $P$ : for an input  $I$  over  $\text{sch}(P)$ , we first compute the fixpoint  $P_1(I)$ , then the fixpoint  $P_2(P_1(I))$ , etc. The *output of  $P$  on input  $I$* , denoted  $P(I)$ , is defined as  $P_k(P_{k-1}(\dots P_1(I)\dots))$ . It is well known that the output of  $P$  does not depend on the chosen syntactic stratification (if more than one exists). Not all Datalog<sup>-</sup> programs are syntactically stratifiable.

### 3.2.3 Stable model semantics

We now recall the stable model semantics for a Datalog<sup>⊖</sup> program  $P$  (Gelfond and Lifschitz 1988; Sacca and Zaniolo 1990). Let  $I$  be an instance over  $sch(P)$ . Let  $\varphi \in P$ . Let  $V$  be a valuation for  $\varphi$  whose image is contained in  $adom(I) \cup C$ , where  $C$  is the set of all constants appearing in  $P$ . Valuation  $V$  does not have to be satisfying for  $\varphi$  on  $I$ . Together,  $V$  and  $\varphi$  give rise to a ground rule  $\psi$ , obtained from  $\varphi$  by replacing each  $u \in vars(\varphi)$  with  $V(u)$ . We call  $\psi$  a *ground rule of  $\varphi$  with respect to  $I$* . Let  $ground(\varphi, I)$  denote the set of all ground rules of  $\varphi$  with respect to  $I$ . The *ground program of  $P$  on  $I$* , denoted  $ground(P, I)$ , is defined as  $\bigcup_{\varphi \in P} ground(\varphi, I)$ . Note, if  $I = \emptyset$ , the set  $ground(P, I)$  contains only rules whose ground atoms are made with  $C$ , or atoms that are nullary.

Let  $M$  be another instance over  $sch(P)$ . We write  $ground_M(P, I)$  to denote the program obtained from  $ground(P, I)$  as follows:

- (1) remove every rule  $\psi \in ground(P, I)$  for which  $neg_\psi \cap M \neq \emptyset$ ;
- (2) remove the negative (ground) body atoms from all remaining rules.

Note,  $ground_M(P, I)$  is a positive program. We say that  $M$  is a *stable model of  $P$  on input  $I$*  if  $M$  is the output of  $ground_M(P, I)$  on input  $I$ . If so, the semantics of positive Datalog<sup>⊖</sup> programs implies  $I \subseteq M$ . Not all Datalog<sup>⊖</sup> programs have stable models on every input (Gelfond and Lifschitz 1988).

### 3.3 Network and distributed databases

A (*computer*) *network* is a nonempty finite set  $\mathcal{N}$  of *nodes*, which are values in **dom**. Intuitively,  $\mathcal{N}$  represents the identifiers of compute nodes involved in a distributed system. Communication channels (edges) are not explicitly represented because we allow a node  $x$  to send a message to any node  $y$ , as long as  $x$  knows about  $y$  by means of input relations or received messages. For general distributed or cluster computing, the delivery of messages is handled by the network layer, which is abstracted away. But (Datalog) programs can also describe the network layer itself (Loo *et al.* 2009; Hellerstein 2010b), in which case we would restrict attention to programs where nodes only send messages to nodes to which they are explicitly linked; these nodes would again be provided as input.

A *distributed database instance*  $H$  over a network  $\mathcal{N}$  and a database schema  $\mathcal{D}$  is a function that maps every node of  $\mathcal{N}$  to an ordinary finite database instance over  $\mathcal{D}$ . This represents how data over the same schema  $\mathcal{D}$  is spread over a network.

As a small example of a distributed database instance, consider the following instance  $H$  over a network  $\mathcal{N} = \{x, y\}$  and a schema  $\mathcal{D} = \{R/1, S/1\}$ :  $H(x) = \{R(a), S(b)\}$  and  $H(y) = \{R(a), S(c)\}$ . In words: we put facts  $R(a)$  and  $S(b)$  at node  $x$ , and we put facts  $R(a)$  and  $S(c)$  at node  $y$ . Note that it is possible that the same fact is given to multiple nodes.

## 4 Model-based semantics

Here we describe a class of distributed Datalog<sup>⊃</sup> programs that we give a model-based semantics. First, in Section 4.1, we recall the user language Dedalus, that is based on Datalog<sup>⊃</sup> with annotations, in which the programmer can express the functionality of the distributed program. Next, we discuss how to assign a declarative, model-based semantics to Dedalus programs. This semantics consists of applying the stable model semantics to the Dedalus programs after they are transformed into pure Datalog<sup>⊃</sup> programs, i.e., without annotations. We introduce some auxiliary notations and symbols in Section 4.2. Next, in Section 4.3, we give a basic transformation of Dedalus programs in order to apply the stable model semantics. However, this basic transformation has some shortcomings, that we iteratively correct in Sections 4.4 and 4.5.

### 4.1 User language: Dedalus

Our user language for distributed Datalog<sup>⊃</sup> programs is Dedalus (Alvaro et al. 2009; Hellerstein 2010b; Alvaro et al. 2011), here presented as Datalog<sup>⊃</sup> with annotations.<sup>3</sup> Essentially, the language represents updatable memory for the nodes of a network and provides a mechanism for communication between these nodes.

#### 4.1.1 Syntax

Let  $\mathcal{D}$  be a database schema. We write  $\mathbf{B}\{\bar{v}\}$ , where  $\bar{v}$  is a tuple of variables, to denote any sequence  $\beta$  of literals over database schema  $\mathcal{D}$ , such that the variables in  $\beta$  are precisely those in the tuple  $\bar{v}$ . Let  $R(\bar{u})$  denote any atom over  $\mathcal{D}$ . There are three types of Dedalus rules over  $\mathcal{D}$ :

- A *deductive* rule is a normal Datalog<sup>⊃</sup> rule over  $\mathcal{D}$ .
- An *inductive* rule is of the form

$$R(\bar{u}) \bullet \leftarrow \mathbf{B}\{\bar{u}, \bar{v}\}.$$

- An *asynchronous* rule is of the form

$$R(\bar{u}) \mid y \leftarrow \mathbf{B}\{\bar{u}, \bar{v}, y\}.$$

For asynchronous rules, the annotation “ $\mid y$ ” with  $y \in \mathbf{var}$  means that the derived head facts are transferred (“piped”) to the addressee node represented by  $y$ . Deductive, inductive and asynchronous rules will express respectively local computation, updatable memory, and message sending. As in Section 3.2, a Dedalus rule is called *safe* if all its variables occur in at least one positive body atom.

We already provide some intuition of how asynchronous rules operate. There are four conceptual time points involved in the execution of an asynchronous rule: the time when the body is evaluated; the time when the derived fact is sent to the

<sup>3</sup> These annotations correspond to syntactic sugar in the previous presentations of Dedalus.

addressee; the time when the fact arrives at the addressee; and, the time when the arrived fact becomes visible at the addressee. In the model-based semantics presented later, the first two time points coincide and the last two time points coincide; and, there is no upper bound on the interval between these two pairs, although it will be finite.

Now consider the following definition:

#### *Definition 4.1*

A *Dedalus program over a schema  $\mathcal{D}$*  is a set of deductive, inductive and asynchronous Dedalus rules over  $\mathcal{D}$ , such that all rules are safe, and the set of deductive rules is syntactically stratifiable.

In the current work, we will additionally assume that Dedalus programs are constant-free, as is common in the theory of database query languages, and which is not really a limitation, since constants that are important for the program can always be indicated by unary relations in the input.

Let  $\mathcal{P}$  be a Dedalus program. The definitions of  $sch(\mathcal{P})$ ,  $idb(\mathcal{P})$ , and  $edb(\mathcal{P})$  are like for Datalog<sup>+</sup> programs. An *input* for  $\mathcal{P}$  is a *distributed* database instance over some network  $\mathcal{N}$  and the schema  $edb(\mathcal{P})$ .

#### 4.1.2 Semantics sketch

We sketch the main idea behind the semantics of a Dedalus program  $\mathcal{P}$ . We illustrate the semantics in Section 4.1.3.

Let  $H$  be an input distributed database instance for  $\mathcal{P}$ , over a network  $\mathcal{N}$ . The idea is that all nodes  $x \in \mathcal{N}$  run the same program  $\mathcal{P}$  and use their local input fragment  $H(x)$  to do local computation and to send messages. Conceptually, each node of  $\mathcal{N}$  should be thought of as doing local computation steps, indefinitely. During each step, a node reads the following facts: (i) the local input; (ii) some received message facts, generated by asynchronous rules on other nodes or the node itself; and, (iii) the facts derived by inductive rules during the previous step on this same node. Next, the deductive rules are applied to these available facts, to compute a fixpoint  $D$  under the stratified semantics.

Subsequently, the asynchronous and inductive rules are fired in parallel on the deductive fixpoint  $D$ , trying all possible valuations in single-step derivations (i.e., no fixpoint). The asynchronous rules send messages to other nodes or to the same node. Messages arrive after an arbitrary (but finite) delay, where the delay can vary for each message. The inductive rules store facts in the memory of the local node. The effect of an inductive derivation is only visible in the very next step; so, if a fact is to be remembered over multiple steps, it should always be explicitly rederived by inductive rules.

#### 4.1.3 Examples

We consider several examples to demonstrate the three kinds of Dedalus rules, and how they work together. These examples also illustrate the utility of Dedalus when

```

marked(u) | y ← start(u), Node(y).
marked(u)• ← marked(u).
marked(v) ← marked(u), R(u, v).
vert(u) ← R(u, v).
vert(u) ← R(v, u).
missing() ← vert(u), ¬marked(u).
covered() ← ¬missing().

```

Figure 1. Dedalus program for example 1.

applied to some practical problems. Here, we follow the principle that the output on a node  $x$  consists of the facts that are eventually derived during every step of  $x$ .

### Example 1

In this example, we compute reachable vertices on graph data. Consider the Dedalus program  $\mathcal{P}$  in Figure 1. We assume the *edb* relations  $R/2$ ,  $\text{start}/1$ , and  $\text{Node}/1$ . For each node, relation  $R$  describes a local graph, and relation  $\text{start}$  provides certain starting vertices. In any input distributed database instance  $H$  over a network  $\mathcal{N}$ , we assume that for each node, relation  $\text{Node}$  is initialized to contain all nodes of  $\mathcal{N}$ ; intuitively,  $\text{Node}$  can be regarded as an address book for  $\mathcal{N}$ .

Now, the idea is that each node of  $\mathcal{N}$  will check whether all of its local vertices are reachable from the (distributed) start vertices. Communication is needed to share these start vertices, which is accomplished by the asynchronous rule. The receipt of a start vertex initializes a local relation  $\text{marked}/1$  at each node; this relation contains reachable vertices. The inductive rule says that all reachable vertices that we know during the current step, are remembered in the next step. This way, the effect of the communication is preserved. Moreover, the third rule, which is deductive, collects all local graph vertices reachable from the currently known reachable vertices. Note, the inductive rule will cause the result of this deductive computation to be also remembered in the next step, although this effect is not really needed here. The last four rules, which are deductive, check that all local vertices are reachable from the start vertices seen so far; if so, a local flag  $\text{covered}()$  is derived.

In our semantics, we will enforce that all messages eventually arrive. In such a semantics, eventually a node will produce  $\text{covered}()$  during each step iff all its local vertices are reachable from the distributed start vertices.

### Example 2

In this example, we generate a random ordering of a set through asynchronous delivery of messages. Every node generates a random ordering of a local *edb* relation  $S/1$  that represents an input set. We also assume an *edb* relation  $\text{Id}/1$  that contains on each node the identifier of that node; the relation  $\text{Id}$  allows a node to send a message to itself. The idea is that a node sends all elements of  $S$  to itself

```

M(u) | x ← S(u), Id(x).

used(u) ← F(u).
used(u) ← N(u, v).
used(u) ← N(v, u).
new(u) ← M(u), ¬used(u).
eq(u, u) ← S(u).
two() ← new(u), new(v), ¬eq(u, v).
keep(u) ← new(u), ¬two().
notlast(u) ← N(u, v).
last(u) ← F(u), ¬notlast(u).
last(u) ← N(v, u), ¬notlast(u).

started() ← F(u).
F(u)• ← ¬started(), keep(u).
N(u, v)• ← started(), last(u), keep(v).
F(u)• ← F(u).
N(u, v)• ← N(u, v).

```

Figure 2. Dedalus program for example 2.

as messages, and the arbitrary arrival order is used to generate an ordering of the elements. This ordering depends on the execution, and some executions will not lead to orderings if some elements are always jointly delivered.

The corresponding program is shown in Figure 2. We use relation  $M/1$  to send the elements of  $S$ , as accomplished by the single asynchronous rule. The relations  $F/1$  and  $N/2$  represent the ordering of  $S$  so far, and they are considered as the output of the program; the letters “F” and “N” stand for “first” and “next” respectively. For example, a possible ordering of the set  $\{a, b, c, d\}$  could be expressed by the following facts:  $F(d)$ ,  $N(d, c)$ ,  $N(c, b)$ ,  $N(b, a)$ .

Inductive rules are responsible for remembering the iteratively updated versions of  $F$  and  $N$ . The other rules are deductive, and they can conceptually be executed in the order in which they are written. The main technical challenge is to only update the ordering when precisely one element of  $S$  arrives; otherwise, because we have no choice mechanism, we would accidentally give the same ordinal to two different elements. Checking whether we may update the ordering is accomplished through other auxiliary relations. We use a nullary relation `started` as a flag to know whether we still have to initialize relation  $F$  or not.

Note that the program keeps sending all elements of  $S$  through the single asynchronous rule. Alternatively, by adapting the program, we could send the elements only once by making sure the asynchronous rule is fired only once (in parallel for all elements of  $S$ ). In that case, as soon as two elements are later delivered together, the ordering will not contain all elements.

*Example 3*

This example is inspired by commit protocols that were expressed in a precursor language of Dedalus (Alvaro et al. 2009). In particular, we implement a two-phase commit protocol where agents, represented by nodes, vote either “yes” or “no” for transaction identifiers. Such a protocol could be part of a bigger system, where transactions are distributed across agents and each agent may only perform the transaction locally if *all* agents want to do this. A single *coordinator* node is responsible for combining the votes for each transaction identifier  $t$ : the coordinator broadcasts “yes” for  $t$  if all votes for  $t$  are “yes”, and “no” otherwise. Each agent stores the decision of the coordinator.

Because the agents and the coordinator have different roles, we make two separate Dedalus programs.<sup>4</sup> First, the agent nodes are assigned the following simple Dedalus program, whose relations are explained below:

$$\text{vote}(t, x, v) \mid y \leftarrow \text{myVote}(t, v), \text{Id}(x), \text{coord}(y).$$

$$\text{outcome}(t, v) \bullet \leftarrow \text{outcome}(t, v).$$

Here, the *edb* relations are: *myVote*/2 that maps each transaction identifier  $t$  to a local vote “yes” or “no”, *Id*/1 storing the identifier of the agent, and *coord*/1 storing the identifier of the coordinator. Also, the relations *vote*/3 and *outcome*/2 represent respectively the outgoing votes and the final decision by the coordinator.

Second, the coordinator node is assigned the Dedalus program shown in Figure 3. The coordinator has the following *edb* relations: relation *T*/1 containing all transaction identifiers, relations *Y*/1 and *N*/1 containing the constants “yes” and “no” respectively, and relation *agents*/1 containing all voting agents. The coordinator uses an inductive rule to gradually accumulate all votes for each transaction identifier. Votes can have arbitrary delays, but in our model the delays are always finite. In each computation step, the deductive rules at the coordinator recompute a relation *complete* that contains the transaction identifiers for which all votes have been received. When a transaction identifier  $t$  has at least one “no” vote, the coordinator decides “no” for  $t$ , and otherwise the coordinator decides “yes” for  $t$ . The final decision is broadcast to all agents. The coordinator adds the transactions with a decision to a log, so the decision will not be broadcast again.

#### 4.2 Auxiliary notations and relations

Let  $\mathcal{P}$  be a Dedalus program. Let  $R/k \in \text{sch}(\mathcal{P})$ . We will use facts of the form  $R(x, s, a_1, \dots, a_k)$  to express that fact  $R(a_1, \dots, a_k)$  is present at a node  $x$  during its local step  $s$ , with  $s \in \mathbb{N}$ , after the deductive rules are executed. We call  $x$  the *location specifier* and  $s$  the *timestamp*. In order to represent timestamps, we assume  $\mathbb{N} \subseteq \text{dom}$ .

<sup>4</sup> In our formal definitions, all nodes execute the same Dedalus program. However, it is easy to simulate two different programs by giving every node the union of both programs, but using a flag to guard the rules of each program. In this example, we can then assume that one node gets a “coordinator” flag as input, and the other nodes get an “agent” flag as input.

Table 1. Relation names not in  $sch(\mathcal{P})$

Relation names	Meaning
all	Network
time, tsucc, <, ≠	Timestamps
before	Happens-before relation
cand <sub>R</sub> , chosen <sub>R</sub> , other <sub>R</sub> , for each relation name R in $idb(\mathcal{P})$	Messages
hasSender, isSmaller, hasMax, rcvInf	Only a finite number of messages arrive at each step of a node

```

vote(t, x, v)• ← vote(t, x, v).
known(t, x) ← vote(t, x, v).
missing(t) ← T(t), agent(x), ¬known(t, x).
complete(t) ← T(t), ¬missing(t).
decideNo(t) ← votes(t, x, v), N(v).
decideYes(t) ← complete(t), ¬decideNo(t).

outcome(t, v) | y ← decideNo(t), ¬log(t), N(v), agent(y).
outcome(t, v) | y ← decideYes(t), ¬log(t), Y(v), agent(y).
log(t)• ← complete(t).
log(t)• ← log(t).
    
```

Figure 3. Dedalus (coordinator) program for Example 3.

We write  $sch(\mathcal{P})^{LT}$  to denote the database schema obtained from  $sch(\mathcal{P})$  by incrementing the arity of every relation by two. The two extra components will contain the location specifier and timestamp.<sup>5</sup> For an instance  $I$  over  $sch(\mathcal{P})$ ,  $x \in \mathbf{dom}$  and  $s \in \mathbb{N}$ , we write  $I^{\uparrow x, s}$  to denote the facts over  $sch(\mathcal{P})^{LT}$  that are obtained by prepending location specifier  $x$  and timestamp  $s$  to every fact of  $I$ . Also, if  $L$  is a sequence of literals over  $sch(\mathcal{P})$ , and  $x, s \in \mathbf{var}$ , we write  $L^{\uparrow x, s}$  to denote the sequence of literals over  $sch(\mathcal{P})^{LT}$  that is obtained by adding location specifier  $x$  and timestamp  $s$  to the literals in  $L$  (negative literals stay negative).

We also need auxiliary relation names, that are assumed not to be used in  $sch(\mathcal{P})$ ; these are listed in Table 1.<sup>6</sup> The concrete purpose of these relations will become clear in the following subsections.

We define the following schema

$$\mathcal{D}_{\text{time}} = \{\mathbf{time}/1, \mathbf{tsucc}/2, \mathbf{<}/2, \mathbf{\neq}/2\}.$$

<sup>5</sup> The abbreviation ‘LT’ stands for ‘location specifier and timestamp’.

<sup>6</sup> In practice, auxiliary relations can be differentiated from those in  $sch(\mathcal{P})$  by a namespace mechanism.

The relations “<” and “≠” will be written in infix notation in rules. We consider only the following instance over  $\mathcal{D}_{\text{time}}$ :

$$I_{\text{time}} = \{\text{time}(s), \text{tsucc}(s, s + 1) \mid s \in \mathbb{N}\} \\ \cup \{(s < t) \mid s, t \in \mathbb{N} : s < t\} \\ \cup \{(s \neq t) \mid s, t \in \mathbb{N} : s \neq t\}.$$

Intuitively, the instance  $I_{\text{time}}$  provides timestamps together with relations to compare them.

### 4.3 Dynamic choice transformation

Let  $\mathcal{P}$  be a Dedalus program. We describe the *dynamic choice transformation* to transform  $\mathcal{P}$  into a pure Datalog<sup>−</sup> program  $\text{pure}_{\text{ch}}(\mathcal{P})$ . The most technical part of the transformation involves the use of dynamic choice to select an arrival timestamp for each message generated by an asynchronous rule. The actual transformation is presented first; next we give the semantics; and, lastly, we discuss how the transformation can be improved.

#### 4.3.1 Transformation

We incrementally construct  $\text{pure}_{\text{ch}}(\mathcal{P})$ . In particular, for each rule in  $\mathcal{P}$ , we specify what corresponding rule (or rules) should be added to  $\text{pure}_{\text{ch}}(\mathcal{P})$ . For technical convenience, we assume that rules of  $\mathcal{P}$  always contain at least one positive body atom. This assumption allows us to more elegantly enforce that head variables in rules of  $\text{pure}_{\text{ch}}(\mathcal{P})$  also occur in at least one positive body atom.<sup>7</sup> Let  $x, s, t, t' \in \mathbf{var}$  be distinct variables not yet occurring in rules of  $\mathcal{P}$ . We write  $\mathbf{B}\{\bar{v}\}$ , where  $\bar{v}$  is a tuple of variables, to denote any sequence  $\beta$  of literals over  $\text{sch}(\mathcal{P})$ , such that the variables in  $\beta$  are precisely those in  $\bar{v}$ . Also recall the notations and relation names from Section 4.2.

*Deductive rules.* For each deductive rule  $R(\bar{u}) \leftarrow \mathbf{B}\{\bar{u}, \bar{v}\}$  in  $\mathcal{P}$ , we add to  $\text{pure}_{\text{ch}}(\mathcal{P})$  the following rule:

$$R(x, s, \bar{u}) \leftarrow \mathbf{B}\{\bar{u}, \bar{v}\}^{\uparrow x, s}. \quad (1)$$

This rule expresses that deductively derived facts at some node  $x$  during step  $s$  are (immediately) visible within step  $s$  of  $x$ . Note, all atoms in this rule are over  $\text{sch}(\mathcal{P})^{\text{LT}}$ .

*Inductive rules.* For each inductive rule  $R(\bar{u}) \bullet \leftarrow \mathbf{B}\{\bar{u}, \bar{v}\}$  in  $\mathcal{P}$ , we add to  $\text{pure}_{\text{ch}}(\mathcal{P})$  the following rule:

$$R(x, t, \bar{u}) \leftarrow \mathbf{B}\{\bar{u}, \bar{v}\}^{\uparrow x, s}, \text{tsucc}(s, t). \quad (2)$$

This rule expresses that inductively derived facts becomes visible in the *next* step of the *same* node.

<sup>7</sup> This assumption is not really a restriction, since a nullary positive body atom is already sufficient.

*Asynchronous rules.* We use facts of the form  $\text{all}(x)$  to say that  $x$  is a node of the network at hand. We use facts of the form  $\text{cand}_R(x, s, y, t, \bar{a})$  to express that node  $x$  at its step  $s$  sends a message  $R(\bar{a})$  to node  $y$ , and that  $t$  could be the arrival timestamp of this message at  $y$ .<sup>8</sup> Within this context, we use a fact  $\text{chosen}_R(x, s, y, t, \bar{a})$  to say that  $t$  is the *effective* arrival timestamp of this message at  $y$ . Lastly, a fact  $\text{other}_R(x, s, y, t, \bar{a})$  means that  $t$  is *not* the arrival timestamp of the message. Now, for each asynchronous rule

$$R(\bar{u}) \mid y \leftarrow \mathbf{B}\{\bar{u}, \bar{v}, y\},$$

in  $\mathcal{P}$ , letting  $\bar{w}$  be a tuple of new and distinct variables with  $|\bar{w}| = |\bar{u}|$ , we add to  $\text{pure}_{\text{ch}}(\mathcal{P})$  the following rules, for which the intuition is given below:

$$\text{cand}_R(x, s, y, t, \bar{u}) \leftarrow \mathbf{B}\{\bar{u}, \bar{v}, y\}^{\uparrow x, s}, \text{all}(y), \text{time}(t). \quad (3)$$

$$\text{chosen}_R(x, s, y, t, \bar{w}) \leftarrow \text{cand}_R(x, s, y, t, \bar{w}), \neg \text{other}_R(x, s, y, t, \bar{w}). \quad (4)$$

$$\text{other}_R(x, s, y, t, \bar{w}) \leftarrow \text{cand}_R(x, s, y, t, \bar{w}), \text{chosen}_R(x, s, y, t', \bar{w}), t \neq t'. \quad (5)$$

$$R(y, t, \bar{w}) \leftarrow \text{chosen}_R(x, s, y, t, \bar{w}). \quad (6)$$

Rule (3) represents the messages that are sent. It evaluates the body of the original asynchronous rule, verifies that the addressee is within the network by using relation  $\text{all}$ , and it generates all possible candidate arrival timestamps.

Now remains the matter of actually choosing *one* arrival timestamp amongst all these candidates. Intuitively, rule (4) selects an arrival timestamp for a message with the condition that this timestamp is not yet ignored, as expressed with relation  $\text{other}_R$ . Also, looking at rule (5), a possible arrival timestamp  $t$  becomes ignored if there is already a chosen arrival timestamp  $t'$  with  $t \neq t'$ . Together, both rules have the effect that exactly one arrival timestamp will be chosen under the stable model semantics. This technical construction is due to Saccà and Zaniolo (1990), who show how to express dynamic choice under the stable model semantics.

Rule (6) represents the actual arrival of an  $R$ -message with the chosen arrival timestamp: the data-tuple in the message becomes part of the addressee's state for relation  $R$ . When the addressee reads relation  $R$ , it thus transparently reads the arrived  $R$ -messages.

Note, if multiple asynchronous rules in  $\mathcal{P}$  have the same head predicate  $R$ , only new  $\text{cand}_R$ -rules have to be added because the rules (4)–(6) are general for all  $R$ -messages.

Note that if there are asynchronous rules in  $\mathcal{P}$ , program  $\text{pure}_{\text{ch}}(\mathcal{P})$  is not syntactically stratifiable if a  $\text{cand}_R$ -rule contains a body atom that (indirectly) negatively depends on  $R$ .<sup>9</sup> In that case,  $\text{pure}_{\text{ch}}(\mathcal{P})$  might not even be locally stratifiable (Apt and Bol 1994).

<sup>8</sup> Here, “cand” abbreviates “candidate”.

<sup>9</sup> Indeed,  $\text{cand}_R$  is used to compute  $R$ , but  $R$  is also used to compute  $\text{cand}_R$ , giving a cycle through negation.

## 4.3.2 Semantics

Now we define the semantics of  $pure_{ch}(\mathcal{P})$ . Let  $H$  be an input distributed database instance for  $\mathcal{P}$ , over a network  $\mathcal{N}$ . Using the notations from Section 4.2, we define  $decl(H)$  to be the following database instance over the schema  $edb(\mathcal{P})^{LT} \cup \{a11/1\} \cup \mathcal{D}_{time}$ :

$$decl(H) = \{R(x, s, \bar{a}) \mid x \in \mathcal{N}, s \in \mathbb{N}, R(\bar{a}) \in H(x)\} \\ \cup \{a11(x) \mid x \in \mathcal{N}\} \cup I_{time}.$$

In words: we make for each node its input facts available at all timestamps; we provide the set of all nodes; and,  $I_{time}$  provides the timestamps with comparison relations.<sup>10</sup> Note, instance  $decl(H)$  is infinite because  $\mathbb{N}$  is infinite.

The stable model semantics for Datalog<sup>-</sup> programs is reviewed in Section 3.2.3. Consider now the following definition:

*Definition 4.2*

For an input distributed database instance  $H$  for  $\mathcal{P}$ , we call any stable model of  $pure_{ch}(\mathcal{P})$  on input  $decl(H)$  a *choice-model* of  $\mathcal{P}$  on input  $H$ .

## 4.3.3 Possible improvement

We illustrate a shortcoming of the dynamic choice transformation. Consider the Dedalus program  $\mathcal{P}$  in Figure 4. We assume that in each input distributed database, the  $edb$  relation  $Id/1$  contains on each node just the identifier of this node. This way, the node can send messages to itself. Relation  $T$  is the intended output relation of  $\mathcal{P}$ . The idea is that a node sends  $A()$  to itself continuously. When  $A()$  arrives, we send  $B()$ , but we also want to create an output fact  $T()$ . We only create  $T()$  when  $B()$  is absent. When  $B()$  is received, it is remembered by inductive rules. Now, we see that the delivery of at least one  $A()$  is necessary to cause a  $B()$  to be sent. This creates the expectation that  $T()$  is always created: at least one  $A()$  is delivered before any  $B()$ . This intuition can be formalized as *causality* (Attiya and Welch 2004) (see also Section 5.2.1).

However, this intuition is violated by some choice-models of  $\mathcal{P}$ , as we demonstrate next. Consider the input distributed database instance  $H$  over a singleton network  $\{z\}$  that assigns the fact  $Id(z)$  to  $z$ . Now, consider the following choice-model  $M$  of  $\mathcal{P}$  on  $H$ :<sup>11</sup>

$$M = decl(H) \cup M_A^{snd} \cup M_A^{rcv} \cup M_B^{snd} \cup M_B^{rcv},$$

<sup>10</sup> For simplicity we already include relation  $<$  in this definition, although this relation will only be used later.

<sup>11</sup> Using straightforward arguments, it can indeed be shown that  $M$  is a stable model of  $pure_{ch}(\mathcal{P})$  on  $decl(H)$ .

$$\begin{array}{l}
 A() \mid \mathbf{x} \leftarrow \text{Id}(\mathbf{x}) \cdot \\
 B() \mid \mathbf{x} \leftarrow A(), \text{Id}(\mathbf{x}) \cdot \\
 \\
 T() \leftarrow A(), \neg B() \cdot \\
 T() \bullet \leftarrow T() \cdot \\
 B() \bullet \leftarrow B() \cdot
 \end{array}$$

Figure 4. Dedalus program sensitive to non-causality.

where

$$\begin{aligned}
 M_A^{\text{snd}} &= \{ \text{cand}_A(z, s, z, t) \mid s, t \in \mathbb{N} \} \\
 &\quad \cup \{ \text{chosen}_A(z, s, z, s + 1) \mid s \in \mathbb{N} \} \\
 &\quad \cup \{ \text{other}_A(z, s, z, t) \mid s, t \in \mathbb{N}, t \neq s + 1 \}; \\
 \\
 M_A^{\text{rcv}} &= \{ A(z, s) \mid s \in \mathbb{N}, s \geq 1 \}; \\
 \\
 M_B^{\text{snd}} &= \{ \text{cand}_B(z, s, z, t) \mid s, t \in \mathbb{N}, s \geq 1 \} \\
 &\quad \cup \{ \text{chosen}_B(z, 1, z, 0) \} \\
 &\quad \cup \{ \text{chosen}_B(z, s, z, s + 1) \mid s \in \mathbb{N}, s \geq 2 \} \\
 &\quad \cup \{ \text{other}_B(z, 1, z, t) \mid t \in \mathbb{N}, t \neq 0 \} \\
 &\quad \cup \{ \text{other}_B(z, s, z, t) \mid s, t \in \mathbb{N}, s \geq 2, t \neq s + 1 \}; \\
 \\
 M_B^{\text{rcv}} &= \{ B(z, s) \mid s \in \mathbb{N} \}.
 \end{aligned}$$

In  $M_B^{\text{snd}}$ , note that one  $B$ -message is sent at timestamp 1 of  $z$ , and arrives at timestamp 0 of  $z$ . We immediately see that this message is peculiar: we should not be able to send a message to arrive in the past. Because of the stray message  $B()$ , the fact  $B()$  exists at all timestamps: it arrives at timestamp 0 and is henceforth persisted by the inductive rule for relation  $B$ ; this is modeled by set  $M_B^{\text{rcv}}$ . Subsequently, there are no ground rules of the form  $T(z, s) \leftarrow A(z, s)$  with  $s \in \mathbb{N}$  in the ground program  $\text{ground}_M(C, I)$ , where  $C = \text{pure}_{\text{ch}}(\mathcal{P})$  and  $I = \text{decl}(H)$ .

In the next subsection, we exclude such unintuitive stable models using an extended transformation of Dedalus programs.

### 4.4 Causality transformation

Let  $\mathcal{P}$  be a Dedalus program. In this section, we present the causality transformation  $\text{pure}_{\text{ca}}(\mathcal{P})$  that extends  $\text{pure}_{\text{ch}}(\mathcal{P})$  to exclude the unintuitive stable models that we have encountered in the previous subsection. We first present the new transformation, and then we discuss how the transformation can still be improved.

#### 4.4.1 Transformation

We define  $\text{pure}_{\text{ca}}(\mathcal{P})$  again incrementally. First, we transform deductive and inductive rules just as in  $\text{pure}_{\text{ch}}(\mathcal{P})$ .

Next, we use facts of the form  $\text{before}(x, s, y, t)$  to express that local step  $s$  of node  $x$  happens before local step  $t$  of node  $y$ . Regardless of  $\mathcal{P}$ , we always add the following rules to  $\text{pure}_{ca}(\mathcal{P})$ :

$$\text{before}(x, s, x, t) \leftarrow \text{all}(x), \text{tsucc}(s, t). \quad (7)$$

$$\text{before}(x, s, y, t) \leftarrow \text{before}(x, s, z, u), \text{before}(z, u, y, t). \quad (8)$$

Rule (7) expresses that on every node, a step happens before the next step. Rule (8) makes relation  $\text{before}$  transitive.

Now, for each asynchronous rule

$$R(\bar{u}) \mid y \leftarrow \mathbf{B}\{\bar{u}, \bar{v}, y\},$$

in  $\mathcal{P}$ , we add to  $\text{pure}_{ca}(\mathcal{P})$  the previous transformation rules (4), (5) and (6) (omitting the  $\text{cand}_R$ -rule), and we add the following new rules, where  $\bar{w}$  is a tuple of new and distinct variables with  $|\bar{w}| = |\bar{u}|$ , and  $x, s$ , and  $t$  are also new variables:

$$\begin{aligned} \text{cand}_R(x, s, y, t, \bar{u}) \leftarrow \mathbf{B}\{\bar{u}, \bar{v}, y\}^{\uparrow x, s}, \text{all}(y), \text{time}(t), \\ \neg \text{before}(y, t, x, s). \end{aligned} \quad (9)$$

$$\text{before}(x, s, y, t) \leftarrow \text{chosen}_R(x, s, y, t, \bar{w}). \quad (10)$$

Like the old rule (3), rule (9) represents the messages that are sent, but now candidate arrival timestamps are restricted by relation  $\text{before}$  to enforce causality. Intuitively, this restriction prevents cycles from occurring in relation  $\text{before}$ . This aligns with the semantics of a real distributed system, where the happens-before relation is a strict partial order (Attiya and Welch 2004) (see also Section 5.2.1).

Rule (10) adds the causal restriction that the local step of the sender happens before the arrival step of the addressee. Together with the previously introduced rules (7) and (8), this will make sure that when the addressee later *causally* replies to the sender, the reply — as generated by a rule of the form (9) — will arrive after this first send-step of the sender.

#### Remark 1

The new program  $\text{pure}_{ca}(\mathcal{P})$  excludes unintuitive models like the one in Section 4.3.3. In the context of that particular example, it will be impossible to exhibit a stable model of  $\text{pure}_{ca}(\mathcal{P})$  in which  $B()$  is sent to timestamp 0. Indeed,  $B()$  can only be sent starting from timestamp 1; timestamp 0 at  $z$  (locally) happens before timestamp 1 at  $z$ ; and, the negative  $\text{before}$ -literal in rule (9) will prevent sending from timestamp 1 at  $z$  to timestamp 0 at  $z$ . Also in scenarios where different nodes  $x$  and  $y$  send messages to each other, when node  $x$  replies to a message of node  $y$  sent at timestamp  $s$  of  $y$ , node  $x$  cannot send the reply to a timestamp  $t$  of  $y$  with  $t < s$ .

#### 4.4.2 Semantics

The semantics of the causality transformation is the same as for the dynamic choice transformation:

*Definition 4.3*

For an input distributed database instance  $H$  for  $\mathcal{P}$ , we call any stable model of  $pure_{ca}(\mathcal{P})$  on input  $decl(H)$  a *causal model* of  $\mathcal{P}$  on input  $H$ .

4.4.3 Possible improvement

We illustrate a shortcoming of the causality transformation. Consider the Dedalus program  $\mathcal{P}$  in Figure 5. We assume that in each input distributed database, the *edb* relation `contact/1` contains intended recipients of messages. Relation  $T$  serves as the output relation of  $\mathcal{P}$ . The idea is that a node sends  $A()$  to its recipients continuously. When  $A()$  arrives, a recipient sets a local flag `first()`. Later, when a second  $A()$  arrives, the recipient creates an output fact  $T()$  that we remember by means of inductive rules. Intuitively, we expect that  $T()$  is always created because the fact  $A()$  is sent infinitely often to a recipient, making this recipient witness the arrival of  $A()$  at (hopefully) two distinct moments.

However, this intuition is violated by some causal models of  $\mathcal{P}$ . Consider the input distributed database instance  $H$  over a network  $\{x, y\}$  that (only) assigns the fact `contact(y)` to  $x$ . Now, consider the following causal model  $M$  of  $\mathcal{P}$  on  $H$ :<sup>12</sup>

$$M = decl(H) \cup M_A^{snd} \cup M_A^{rcv} \cup M^{before},$$

where

$$M_A^{snd} = \{cand_A(x, s, y, t) \mid s, t \in \mathbb{N}\} \\ \cup \{chosen_A(x, s, y, 0) \mid s \in \mathbb{N}\} \\ \cup \{other_A(x, s, y, t) \mid s, t \in \mathbb{N}, t \neq 0\};$$

$$M_A^{rcv} = \{A(y, 0)\} \\ \cup \{first(y, s) \mid s \in \mathbb{N}, s \geq 1\};$$

$$M^{before} = \{before(x, s, x, t) \mid s, t \in \mathbb{N}, s < t\}; \\ \cup \{before(y, s, y, t) \mid s, t \in \mathbb{N}, s < t\}; \\ \cup \{before(x, s, y, t) \mid s, t \in \mathbb{N}\}.$$

In this causal model, all instances of message  $A()$  that  $x$  sends to  $y$  arrive at timestamp 0 of  $y$ . For this reason, node  $y$  cannot witness two different arrivals of message  $A()$ . In practice, however, node  $y$  cannot receive an *infinite* number of messages during a timestamp, and the deliveries of the  $A()$  messages would be spread out more evenly in time. So, in the next subsection, we will additionally exclude such infinite message arrivals, to obtain our final transformation of Dedalus programs.

4.5 Causality-finiteness transformation

Let  $\mathcal{P}$  be a Dedalus program. As seen in the previous subsection, program  $pure_{ca}(\mathcal{P})$  allows an infinite number of messages to arrive at any step of a node. This does

<sup>12</sup> Using straightforward arguments, it can be shown that  $M$  is a stable model of  $pure_{ca}(\mathcal{P})$  on  $decl(H)$ .

$A() \mid y \leftarrow \text{contact}(y).$ $\text{first}() \bullet \leftarrow A().$ $\text{first}() \bullet \leftarrow \text{first}().$ $T() \leftarrow \text{first}(), A().$ $T() \bullet \leftarrow T().$
---

Figure 5. Dedalus program sensitive to infinite message grouping.

not happen in any real-world distributed system; indeed, no node has to process an infinite number of messages at any given moment. We consider this to be an additional restriction that must be explicitly enforced. To this purpose, we present in this section the causality-finiteness transformation  $\text{pure}(\mathcal{P})$  that extends  $\text{pure}_{\text{ca}}(\mathcal{P})$ .

We will approach this problem as follows. Suppose there are an infinite number of messages that arrive at some node  $y$  during its step  $t$ . Since in a network there are only a finite number of nodes and a node can only send a finite number of messages during each step (the input domain is finite), there must be at least one node  $x$  that sends messages to step  $t$  of  $y$  during an infinite number of steps of  $x$ . Hence, there is no maximum value amongst the corresponding send-timestamps of  $x$ . Thus, in order to prevent the arrival of an infinite number of messages at step  $t$  of  $y$ , it will be sufficient to demand that there always *is* such a maximum send-timestamp for every sender. Below, we will implement this strategy with some concrete rules in  $\text{pure}(\mathcal{P})$ .

#### 4.5.1 Transformation

We define  $\text{pure}(\mathcal{P})$  as  $\text{pure}_{\text{ca}}(\mathcal{P})$  extended as follows. The additional rules can be thought of as being relative to an addressee and a step of this addressee, represented by the variables  $y$  and  $t$  respectively.

We use a fact  $\text{rcvInf}(y, t)$  to express that node  $y$  receives an infinite number of messages during its step  $t$ . First, we add the following rule to  $\text{pure}(\mathcal{P})$  for *each* relation  $\text{chosen}_R$  that results from the transformation of asynchronous rules in  $\text{pure}_{\text{ca}}(\mathcal{P})$ , where  $x, s, y,$  and  $t$  are variables and  $\bar{w}$  is a tuple of distinct variables disjoint from the previous ones with  $|\bar{w}|$  the arity of relation  $R$  in  $\text{sch}(\mathcal{P})$ :

$$\text{hasSender}(y, t, x, s) \leftarrow \text{chosen}_R(x, s, y, t, \bar{w}), \neg \text{rcvInf}(y, t). \quad (11)$$

This rule intuitively means that as long as addressee  $y$  has not received an infinite number of messages during its step  $t$ , we register the senders and their send-timestamps.

Recall the auxiliary relations defined in Section 4.2. Next, we add to  $\text{pure}(\mathcal{P})$  the following rules, for which the intuition is provided below:

$$\text{isSmaller}(y, t, x, s) \leftarrow \text{hasSender}(y, t, x, s), \text{hasSender}(y, t, x, s'), \quad (12)$$

$$s < s'.$$

$$\text{hasMax}(y, t, x) \leftarrow \text{hasSender}(y, t, x, s), \neg \text{isSmaller}(y, t, x, s). \quad (13)$$

$$\text{rcvInf}(y, t) \leftarrow \text{hasSender}(y, t, x, s), \neg \text{hasMax}(y, t, x). \quad (14)$$

Rule (12) checks for each sender and each of its send-timestamps whether there is a later send-timestamp of that same sender. Rule (13) tries to find a maximum send-timestamp. Finally, rule (14) derives a `rcvInf`-fact if no maximum send-timestamp was found for at least one sender.

We will show in Section 5.3.1 that in any stable model, the above rules make sure that every node receives only a finite number of messages at every step.

#### 4.5.2 Semantics

The semantics of the causality-finiteness transformation is again the same as for the dynamic choice transformation and the causality transformation:

##### Definition 4.4

For an input distributed database instance  $H$  for  $\mathcal{P}$ , we call any stable model of  $\text{pure}(\mathcal{P})$  on input  $\text{decl}(H)$  a *causal-finite model* of  $\mathcal{P}$  on input  $H$ .

We will refer to a causal-finite model also simply as *model*.

## 5 Correctness

In Section 4, we have described the computation of a distributed  $\text{Datalog}^-$  program by means of stable models. By using suitable rules, we have excluded some unintuitive stable models. But at this point, we are still not sure whether the remaining stable models really correspond to the execution of a distributed system. We fill that gap in this section: we show that each remaining stable model corresponds to an execution of the distributed  $\text{Datalog}^-$  program under an operational semantics, and vice versa. We call such an execution a *run*, and we will only be concerned with so-called *fair* runs, where each node is made active infinitely often and all sent messages are eventually delivered.

We extract from each run  $\mathcal{R}$  a *trace*, denoted  $\text{trace}(\mathcal{R})$ , which is a set of facts that shows in detail what each node computes during each step. We will make this concrete in the following subsections. But we can already state our main result, as follows:

##### Theorem 4

Let  $\mathcal{P}$  be a Dedalus program. For each input distributed database instance  $H$  for  $\mathcal{P}$ ,

- (i) for every fair run  $\mathcal{R}$  of  $\mathcal{P}$ , there is a model  $M$  of  $\mathcal{P}$  such that  $\text{trace}(\mathcal{R}) = M|_{\text{sch}(\mathcal{P})^{\text{LT}}}$ , and
- (ii) for every model  $M$  of  $\mathcal{P}$ , there is a fair run  $\mathcal{R}$  of  $\mathcal{P}$  such that  $\text{trace}(\mathcal{R}) = M|_{\text{sch}(\mathcal{P})^{\text{LT}}}$ .

First, Section 5.1 formalizes runs and traces of runs. The proof of item (i) of the theorem is described in Section 5.2. The proof of item (ii), which is the most difficult, is described in Section 5.3. We only describe the crucial reasoning steps of

the proofs; the intricate technical details can be found in the online appendix to the paper.

### 5.1 Operational semantics

In this section, we give an operational semantics for Dedalus that is in line with earlier formal work on declarative networking (Deutsch et al. 2006; Navarro and Rybalchenko 2009; Grumbach and Wang 2010; Ameloot et al. 2011; Abiteboul et al. 2011).

Let  $\mathcal{P}$  be a Dedalus program, and let  $H$  be an input distributed database instance for  $\mathcal{P}$ , over a network  $\mathcal{N}$ . The essence of the operational semantics is as follows. Every node of  $\mathcal{N}$  runs program  $\mathcal{P}$ , and a node has access only to its own local state and any received messages. The nodes are made active one by one in some arbitrary order, and this continues an infinite number of times. During each active moment of a node  $x$ , called a *local (computation) step*, node  $x$  receives message facts and applies its deductive, inductive and asynchronous rules. Concretely, the deductive rules, forming a stratified Datalog<sup>-</sup> subprogram, are applied to the incoming messages and the previous state of  $x$ . Next, the inductive rules are applied to the output of the deductive subprogram, and these allow  $x$  to store facts in its memory: these facts become visible in the next local step of  $x$ . Finally, the asynchronous rules are also applied to the output of the deductive subprogram, and these allow  $x$  to send facts to the other nodes or to itself. These facts become visible at the addressee after some arbitrary delay, which represents asynchronous communication, as occurs for instance on the Internet. We assume that all messages are eventually delivered (and are thus never lost). We will refer to local steps simply as “steps”.

We make the above sketch more concrete in the next subsections.

#### 5.1.1 Configurations

Let  $\mathcal{P}$ ,  $H$ , and  $\mathcal{N}$  be as above. A configuration describes the network at a certain point in its evolution. Formally, a *configuration* of  $\mathcal{P}$  on  $H$  is a pair  $\rho = (st, bf)$  where

- $st$  is a function mapping each node of  $\mathcal{N}$  to an instance over  $sch(\mathcal{P})$ ; and,
- $bf$  is a function mapping each node of  $\mathcal{N}$  to a set of pairs of the form  $(i, \mathbf{f})$ , where  $i \in \mathbb{N}$  and  $\mathbf{f}$  is a fact over  $idb(\mathcal{P})$ .

We call  $st$  and  $bf$  the *state* and (*message*) *buffer* respectively. The state says for each node what facts it has stored in its memory, and the message buffer  $bf$  says for each node what messages have been sent to it but that are not yet received. The reason for having numbers  $i$ , called *send-tags*, attached to facts in the image of  $bf$  is merely a technical convenience: these numbers help separate multiple instances of the same fact when it is sent at different moments (to the same addressee), and these send-tags will not be visible to the Dedalus program. For example, if the buffer of a node  $x$  simultaneously contains pairs  $(3, \mathbf{f})$  and  $(7, \mathbf{f})$ , this means that  $\mathbf{f}$  was sent to  $x$  during the operational network transitions with indices 3 and 7, and that both

particular instances of  $f$  are not yet delivered to  $x$ . This will become more concrete in Section 5.1.3.

The *start configuration of  $\mathcal{P}$  on input  $H$* , denoted  $start(\mathcal{P}, H)$ , is the configuration  $\rho = (st, bf)$  defined by  $st(x) = H(x)$  and  $bf(x) = \emptyset$  for each  $x \in \mathcal{N}$ . In words: for every node, the state is initialized with its local input fragment in  $H$ , and there are no sent messages.

### 5.1.2 Subprograms

We look at the operations that are executed locally during each step of a node. We have mentioned that the three types of Dedalus rules each have their own purpose in the operational semantics. For this reason, we split the program  $\mathcal{P}$  into three subprograms, that contain respectively the deductive, inductive and asynchronous rules. In Section 5.1.3, we describe how these subprograms are used in the operational semantics.

- First, we define  $deduc_{\mathcal{P}}$  to be the Datalog<sup>⊖</sup> program consisting of precisely all deductive rules of  $\mathcal{P}$ .
- Secondly, we define  $induc_{\mathcal{P}}$  to be the Datalog<sup>⊖</sup> program consisting of all inductive rules of  $\mathcal{P}$  after the annotation ‘•’ in their head is removed.
- Thirdly, we define  $async_{\mathcal{P}}$  to be the Datalog<sup>⊖</sup> program consisting of precisely all rules

$$T(y, \bar{u}) \leftarrow \mathbf{B}\{\bar{u}, y\}$$

where

$$T(\bar{u}) \mid y \leftarrow \mathbf{B}\{\bar{u}, y\}$$

is an asynchronous rule of  $\mathcal{P}$ . So, we basically put the variable  $y$  as the first component in the (extended) head atom. The intuition for the generated head facts is that the first component will represent the addressee.

Note that the programs  $deduc_{\mathcal{P}}$ ,  $induc_{\mathcal{P}}$  and  $async_{\mathcal{P}}$  are just Datalog<sup>⊖</sup> programs over the schema  $sch(\mathcal{P})$ , or a subschema thereof. Moreover,  $deduc_{\mathcal{P}}$  is syntactically stratifiable because the deductive rules in every Dedalus program must be syntactically stratifiable. It is possible however that  $induc_{\mathcal{P}}$  and  $async_{\mathcal{P}}$  are not syntactically stratifiable. Now we define the semantics of each of these three subprograms.

Let  $I$  be a database instance over  $sch(\mathcal{P})$ . During each step of a node, the intuition of the deductive rules is that they “complete” the available facts by adding all new facts that can be logically derived from them. This calls for a fixpoint semantics, and for this reason, we define the *output of  $deduc_{\mathcal{P}}$  on input  $I$* , denoted as  $deduc_{\mathcal{P}}(I)$ , to be given by the stratified semantics. This implies  $I \subseteq deduc_{\mathcal{P}}(I)$ . Importantly,  $I$  is allowed to contain facts over  $idb(\mathcal{P})$ , and the intuition is that these facts were derived during a previous step (by inductive rules) or received as messages (as sent by asynchronous rules). This will become more explicit in Section 5.1.3.

During each step of a node, the intuition behind the inductive rules is that they store facts in the memory of the node, and these stored facts will become visible during the next step. There is no notion of a fixpoint here because facts that will

become visible in the next step are not available in the current step to derive more facts. For this reason, we define the *output of  $\text{induc}_{\mathcal{P}}$  on input  $I$*  to be the set of facts derived by the rules of  $\text{induc}_{\mathcal{P}}$  for all possible satisfying valuations in  $I$ , in just one derivation step. This output is denoted as  $\text{induc}_{\mathcal{P}}(I)$ .

During each step of a node, the intuition behind the asynchronous rules is that they generate message facts that are to be sent around the network. The *output for  $\text{async}_{\mathcal{P}}$  on input  $I$*  is defined in the same way as for  $\text{induc}_{\mathcal{P}}$ , except that we now use the rules of  $\text{async}_{\mathcal{P}}$  instead of  $\text{induc}_{\mathcal{P}}$ . This output is denoted as  $\text{async}_{\mathcal{P}}(I)$ . The intuition for not requiring a fixpoint for  $\text{async}_{\mathcal{P}}$  is that a message fact will arrive at another node, or at a later step of the sender node, and can therefore not be read during sending.

Regarding data complexity (Vardi 1982), for each subprogram the output can be computed in PTIME with respect to the size of its input.

### 5.1.3 Transitions and runs

Transitions formalize how to go from one configuration to another. Here we use the subprograms of  $\mathcal{P}$ . Transitions are chained to form a *run*. Regarding notation, for a set  $m$  of pairs of the form  $(i, \mathbf{f})$ , we define  $\text{untag}(m) = \{\mathbf{f} \mid \exists i \in \mathbb{N} : (i, \mathbf{f}) \in m\}$ .

A *transition with send-tag  $i \in \mathbb{N}$*  is a five-tuple  $(\rho_a, x, m, i, \rho_b)$  such that  $\rho_a = (st_a, bf_a)$  and  $\rho_b = (st_b, bf_b)$  are configurations of  $\mathcal{P}$  on input  $H$ ,  $x \in \mathcal{N}$ ,  $m \subseteq bf_a(x)$ , and, letting

$$\begin{aligned} I &= st_a(x) \cup \text{untag}(m), \\ D &= \text{deduc}_{\mathcal{P}}(I), \\ \delta^{i \rightarrow y} &= \{(i, R(\bar{a})) \mid R(y, \bar{a}) \in \text{async}_{\mathcal{P}}(D)\} \text{ for each } y \in \mathcal{N}, \end{aligned}$$

for  $x$  and each  $y \in \mathcal{N} \setminus \{x\}$  we have

$$\begin{aligned} st_b(x) &= H(x) \cup \text{induc}_{\mathcal{P}}(D), & st_b(y) &= st_a(y), \\ bf_b(x) &= (bf_a(x) \setminus m) \cup \delta^{i \rightarrow x}, & bf_b(y) &= bf_a(y) \cup \delta^{i \rightarrow y}. \end{aligned}$$

We call  $\rho_a$  and  $\rho_b$  respectively the *source* and *target* configuration, and say this transition is *of* the *active* node  $x$ . Intuitively, the transition expresses that  $x$  reads its old state together with the received facts in  $\text{untag}(m)$  (thus without the tags), and describes the subsequent computation: subprogram  $\text{deduc}_{\mathcal{P}}$  completes the available information; the new state of  $x$  consists of the input facts of  $x$  united with all facts derived by subprogram  $\text{induc}_{\mathcal{P}}$ ; and, subprogram  $\text{async}_{\mathcal{P}}$  generates *messages*, whose first component indicates the addressee.<sup>13</sup> Note,  $\text{induc}_{\mathcal{P}}$  and  $\text{async}_{\mathcal{P}}$  do not influence each other, and can be thought of as being executed in parallel. Also, for each  $y \in \mathcal{N}$ , the set  $\delta^{i \rightarrow y}$  contains all messages addressed to  $y$ , with send-tag  $i$  attached. Messages with an addressee outside the network are ignored. This way of defining local computation closely corresponds to that of the language Webdamlog (Abiteboul et al. 2011). If  $m = \emptyset$ , we call the transition a *heartbeat*.

<sup>13</sup> Note, input facts are preserved by the transition. This aligns with the design of Dedalus, where we do not allow facts to be retracted; only negation as failure is permitted.

A run  $\mathcal{R}$  of  $\mathcal{P}$  on input  $H$  is an infinite sequence of transitions, such that (i) the source configuration of the first transition is  $start(\mathcal{P}, H)$ , (ii) the target configuration of each transition is the source configuration of the next transition, and (iii) the transition at ordinal  $i$  of the sequence uses send-tag  $i$ . Ordinals start at 0 for technical convenience. The resulting transition system is highly non-deterministic because in each transition we can choose the active node and also what messages to deliver; the latter choice is represented by the set  $m$  from above.

*Remark 2 (Parallel transitions)*

Transitions as defined here can simulate *parallel* transitions in which multiple nodes are active at the same time and receive messages from their respective buffers. Indeed, if we would have multiple nodes active during a parallel transition, they would receive messages from their buffers in isolation, and this can be represented by a chain of transitions in which these nodes receive one after the other precisely the messages that they received in the parallel transition. For this reason, we limit our attention to transitions with single active nodes.

#### 5.1.4 Fairness and arrival function

In the literature on process models it is customary to require certain fairness conditions on the execution of a system, for instance to exclude some extreme situations that are expected not to happen in reality (Francez 1986; Apt *et al.* 1988; Lamport 2000b).

Let  $\mathcal{R}$  be a run of  $\mathcal{P}$  on  $H$ . For every transition  $i \in \mathbb{N}$ , let  $\rho_i = (st_i, bf_i)$  denote the source configuration of transition  $i$ . Now,  $\mathcal{R}$  is called *fair* if:

- every node is the active node in an infinite number of transitions of  $\mathcal{R}$ ; and,
- for every transition  $i \in \mathbb{N}$ , for every  $y \in \mathcal{N}$ , for every pair  $(j, \mathbf{f}) \in bf_i(y)$ , there is a transition  $k$  with  $i \leq k$  in which  $(j, \mathbf{f})$  is delivered to  $y$ .

Intuitively, the fairness conditions disallow starvation: every node does an infinite number of local computation steps and every sent message is eventually delivered. We consider only fair runs in this paper. Note, a fair run exists for every input because heartbeats remain possible even when there are no messages to deliver.

In the second condition about message deliveries, it is possible that  $k = i$ , and in that case  $(j, \mathbf{f})$  is delivered in the transition immediately following configuration  $\rho_i$ . Because the pair  $(j, \mathbf{f})$  can be in the message buffer of multiple nodes, this  $k$  is not unique for the pair  $(j, \mathbf{f})$  by itself. But, when we also consider the addressee  $y$ , it follows from the operational semantics that this  $k$  is unique for the triple  $(j, y, \mathbf{f})$ .

This reasoning gives rise to a function  $\alpha_{\mathcal{R}}$ , called the *arrival function for  $\mathcal{R}$* , that is defined as follows: for every transition  $i$ , for every node  $y$ , for every message  $\mathbf{f}$  sent to addressee  $y$  during  $i$ , the function  $\alpha_{\mathcal{R}}$  maps  $(i, y, \mathbf{f})$  to the transition ordinal  $k$  in which  $(i, \mathbf{f})$  is delivered to  $y$ . We always have  $\alpha_{\mathcal{R}}(i, y, \mathbf{f}) > i$ . Indeed, the delivery of a message can only happen after it was sent. So, when the delivery of one message causes another to be sent, then the second one is delivered in a later transition. This is related to the topic of causality that we have introduced in Section 4. This topic will also be further discussed in Sections 5.2 and 5.3.

### 5.1.5 Timestamps and trace

For each transition  $i$  of a run, we define the *timestamp* of the active node  $x$  during  $i$  to be the number of transitions of  $x$  that come strictly before  $i$ . This can be thought of as the *local* (zero-based) clock of  $x$  during  $i$ , and is denoted  $loc_{\mathcal{R}}(i)$ . For example, suppose we have the following sequence of active nodes:  $x, y, y, x, x$ , etc. If we would write the timestamps next to the nodes, we get this sequence:  $(x, 0), (y, 0), (y, 1), (x, 1), (x, 2)$ , etc.

As a counterpart to function  $loc_{\mathcal{R}}(\cdot)$ , for each  $(x, s) \in \mathcal{N} \times \mathbb{N}$  we define  $glob_{\mathcal{R}}(x, s)$  to be the transition ordinal  $i$  of  $\mathcal{R}$  such that  $x$  is the active node in transition  $i$  and  $loc_{\mathcal{R}}(i) = s$ . In words: we find the transition in which node  $x$  does its local computation step with timestamp  $s$ . It follows from the definition of  $loc_{\mathcal{R}}(\cdot)$  that  $glob_{\mathcal{R}}(x, s)$  is uniquely defined.

Let  $\mathcal{R}$  be a run of  $\mathcal{P}$  on input  $H$ . Recall that  $H$  is over network  $\mathcal{N}$ . We now capture the computed data during  $\mathcal{R}$  as a set of facts that we call the *trace*. For each transition  $i \in \mathbb{N}$ , let  $x_i$  denote the active node, and let  $D_i$  denote the output of subprogram  $deduc_{\mathcal{P}}$  during  $i$ . The operational semantics implies that  $D_i$  consists of (i) the input *edb*-facts at  $x_i$ ; (ii) the inductively derived facts during the previous step of  $x_i$  (if  $loc_{\mathcal{R}}(i) \geq 1$ ); (iii) the messages delivered during transition  $i$ ; and, (iv) all facts deductively derived from the previous ones. So, intuitively,  $D_i$  contains *all* local facts over  $sch(\mathcal{P})$  that  $x_i$  has during transition  $i$ .

Recall the notations of Section 4.2. Now, the trace of  $\mathcal{R}$  is the following instance over  $sch(\mathcal{P})^{LT}$ :

$$trace(\mathcal{R}) = \bigcup_{i \in \mathbb{N}} D_i^{\uparrow x_i, loc_{\mathcal{R}}(i)}.$$

The trace shows in detail what happens in the run, in terms of what facts are available on the nodes during which of their steps.

## 5.2 Run to model

Let  $\mathcal{P}$  be a Dedalus program and let  $H$  be an input distributed database instance for  $\mathcal{P}$ , over a network  $\mathcal{N}$ . Let  $\mathcal{R}$  be a fair run of  $\mathcal{P}$  on input  $H$ . We show there is a model  $M$  of  $\mathcal{P}$  on  $H$  such that  $trace(\mathcal{R}) = M|_{sch(\mathcal{P})^{LT}}$ . The main idea is that we translate the transitions of  $\mathcal{R}$  to facts over the schema of  $pure(\mathcal{P})$ .

First, in Section 5.2.1, we extract the happens-before relation on nodes and timestamps from  $\mathcal{R}$ . Next, in Section 5.2.2, we define the desired model  $M$ .

### 5.2.1 Happens-before relation

In the operational semantics, we order the actions of the nodes on a fine-grained global time axis, by ordering the transitions in the runs. By contrast, we now define a partial order on  $\mathcal{N} \times \mathbb{N}$ , saying which steps of nodes must have come before which steps of (other) nodes, without referring to the global ordering imposed by transitions.

First, we extract from  $\mathcal{R}$  the message sending and receiving events. Formally, we define  $mesg(\mathcal{R})$  to be the set of all tuples  $(x, s, y, t, \mathbf{f})$ , with  $\mathbf{f}$  a fact, and denoting

$i = glob_{\mathcal{R}}(x, s)$  and  $j = glob_{\mathcal{R}}(y, t)$ , such that  $\alpha_{\mathcal{R}}(i, y, \mathbf{f}) = j$ , i.e., node  $x$  during step  $s$  sends message  $\mathbf{f}$  to  $y$  that arrives at the step  $t$  of  $y$ , with possibly  $x = y$ . In words:  $msg(\mathcal{R})$  contains the direct relationships between local steps of nodes that arise through message sending.

From  $\mathcal{R}$ , we can now extract the *happens-before* relation (Attiya and Welch 2004) on the set  $\mathcal{N} \times \mathbb{N}$ , which is defined as the smallest relation  $<_{\mathcal{R}}$  on  $\mathcal{N} \times \mathbb{N}$  that satisfies the following three conditions:

- for each  $(x, s) \in \mathcal{N} \times \mathbb{N}$ , we have  $(x, s) <_{\mathcal{R}} (x, s + 1)$ ;
- $(x, s) <_{\mathcal{R}} (y, t)$  whenever for some fact  $\mathbf{f}$  we have  $(x, s, y, t, \mathbf{f}) \in msg(\mathcal{R})$ ;
- $<_{\mathcal{R}}$  is transitive, i.e.,  $(x, s) <_{\mathcal{R}} (z, u) <_{\mathcal{R}} (y, t)$  implies  $(x, s) <_{\mathcal{R}} (y, t)$ .

We call these three cases respectively *local* edges, *message* edges, and *transitive* edges. Naturally, the first two cases express a direct relationship, whereas the third case is more indirect.

Note, if two runs on the same input have the same happens-before relation, they do not necessarily have the same trace. This is because relation  $<_{\mathcal{R}}$  does not talk about the specific messages that arrive at the nodes.

We will now show that  $<_{\mathcal{R}}$  is a strict partial order. Consider first the following property:

*Lemma 1*

For every run  $\mathcal{R}$ , for each  $(x, s) \in \mathcal{N} \times \mathbb{N}$  and  $(y, t) \in \mathcal{N} \times \mathbb{N}$ , if  $(x, s) <_{\mathcal{R}} (y, t)$  then  $glob_{\mathcal{R}}(x, s) < glob_{\mathcal{R}}(y, t)$ .

*Proof*

We can consider a path from  $(x, s)$  to  $(y, t)$  in  $<_{\mathcal{R}}$ . We can substitute each transitive edge in this path with a subpath of non-transitive edges. This results in a path of only non-transitive edges:

$$(x_1, s_1) <_{\mathcal{R}} (x_2, s_2) <_{\mathcal{R}} \dots <_{\mathcal{R}} (x_n, s_n),$$

where  $n \geq 2$ ,  $(x_1, s_1) = (x, s)$  and  $(x_n, s_n) = (y, t)$ . Because there are no transitive edges, for each  $i \in \{1, \dots, n - 1\}$ , the edge  $(x_i, s_i) <_{\mathcal{R}} (x_{i+1}, s_{i+1})$  falls into one of the following two cases:

- $x_i = x_{i+1}$  and  $s_{i+1} = s_i + 1$  (local edge);
- $x_i$  during step  $s_i$  sends a message to  $x_{i+1}$  that arrives in step  $s_{i+1}$  of  $x_{i+1}$  (message edge).

In the first case, it follows from the definition of  $loc_{\mathcal{R}}(\cdot)$  that

$$glob_{\mathcal{R}}(x_i, s_i) < glob_{\mathcal{R}}(x_{i+1}, s_{i+1}).$$

For the second case, by our operational semantics, every message is always delivered in a later transition than the one in which it was sent. So, again we have

$$glob_{\mathcal{R}}(x_i, s_i) < glob_{\mathcal{R}}(x_{i+1}, s_{i+1}).$$

Since this property holds for all the above edges, by transitivity we thus have  $glob_{\mathcal{R}}(x, s) < glob_{\mathcal{R}}(y, t)$ , as desired. □

Corollary 1

For every run  $\mathcal{R}$ , the relation  $<_{\mathcal{R}}$  is a strict partial order on  $\mathcal{N} \times \mathbb{N}$ .

Proof

From its definition, we immediately have that  $<_{\mathcal{R}}$  is transitive. Secondly, irreflexivity for  $<_{\mathcal{R}}$  follows from Lemma 1. □

5.2.2 Definition of  $M$

Now, we define the model  $M$ :

$$M = \text{decl}(H) \cup \bigcup_{i \in \mathbb{N}} \text{trans}_{\mathcal{R}}^{[i]},$$

where  $\text{trans}_{\mathcal{R}}^{[i]}$  for each  $i \in \mathbb{N}$  is an instance over the schema of  $\text{pure}(\mathcal{P})$  that describes transition  $i$  of  $\mathcal{R}$ .<sup>14</sup> Let  $i \in \mathbb{N}$ . We define  $\text{trans}_{\mathcal{R}}^{[i]}$  as

$$\text{trans}_{\mathcal{R}}^{[i]} = \text{caus}_{\mathcal{R}}^{[i]} \cup \text{fin}_{\mathcal{R}}^{[i]} \cup \text{duc}_{\mathcal{R}}^{[i]} \cup \text{snd}_{\mathcal{R}}^{[i]},$$

where each of these sets focuses on different aspects of transition  $i$ , and they are defined next. Regarding notation, let  $<_{\mathcal{R}}$  be the happens-before relation as defined in the preceding subsection; let  $\text{loc}_{\mathcal{R}}(\cdot)$ ,  $\text{glob}_{\mathcal{R}}(\cdot)$ , and  $\alpha_{\mathcal{R}}$  be as defined in Section 5.1; let  $x_i$  denote the active node of transition  $i$ ; and, let us abbreviate  $s_i = \text{loc}_{\mathcal{R}}(i)$ .

*Causality.* We define  $\text{caus}_{\mathcal{R}}^{[i]}$  to consist of all facts  $\text{before}(x, s, x_i, s_i)$  for which  $(x, s) \in \mathcal{N} \times \mathbb{N}$  and  $(x, s) <_{\mathcal{R}} (x_i, s_i)$ . Intuitively,  $\text{caus}_{\mathcal{R}}^{[i]}$  represents the joint result of rules (7), (8), and (10), corresponding to respectively the local edges, transitive edges, and message edges of  $<_{\mathcal{R}}$ .

*Finite Messages.* We define  $\text{fin}_{\mathcal{R}}^{[i]}$  to represent that only a finite number of messages are delivered in transition  $i$ , thus at step  $s_i$  of node  $x_i$ . We proceed as follows. First, let  $\text{senders}_{\mathcal{R}}^{[i]}$  be the set of all pairs  $(x, s) \in \mathcal{N} \times \mathbb{N}$  such that, denoting  $j = \text{glob}_{\mathcal{R}}(x, s)$ , for some fact  $\mathbf{f}$  we have  $\alpha_{\mathcal{R}}(j, x_i, \mathbf{f}) = i$ , i.e., the node  $x$  during its step  $s$  sends a message to  $x_i$  with arrival timestamp  $s_i$ . It follows from the operational semantics that for each  $(x, s) \in \text{senders}_{\mathcal{R}}^{[i]}$  we have  $\text{glob}_{\mathcal{R}}(x, s) < i$ . Now, we define  $\text{fin}_{\mathcal{R}}^{[i]}$  to consist of the following facts:

- the fact  $\text{hasSender}(x_i, s_i, x, s)$  for each  $(x, s) \in \text{senders}_{\mathcal{R}}^{[i]}$ , representing the result of rule (11);
- the fact  $\text{isSmaller}(x_i, s_i, x, s)$  for each  $(x, s) \in \text{senders}_{\mathcal{R}}^{[i]}$  and  $(x, s') \in \text{senders}_{\mathcal{R}}^{[i]}$  with  $s < s'$ , representing the result of rule (12); and,
- the fact  $\text{hasMax}(x_i, s_i, x)$  for each sender-node  $x$  mentioned in  $\text{senders}_{\mathcal{R}}^{[i]}$ , representing the result of rule (13).

We know that in  $\mathcal{R}$  only a finite number of messages arrive at step  $s_i$  of  $x_i$ . Hence, we add no fact  $\text{rcvInf}(x_i, s_i)$  to  $\text{fin}_{\mathcal{R}}^{[i]}$ . This also explains why the specification of the  $\text{hasMax}$ -facts above is relatively simple: there is always a maximum send-timestamp for each sender-node.

<sup>14</sup> Note,  $M$  must include the input  $\text{decl}(H)$  by definition of stable model (see Section 3.2.3).

*Deductive.* Let  $D_i$  denote the output of subprogram  $deduc_{\mathcal{P}}$  during transition  $i$ . We define  $duc_{\mathcal{R}}^{[i]}$  to consist of the facts  $D_i^{\uparrow x_i, s_i}$ . Intuitively,  $duc_{\mathcal{R}}^{[i]}$  represents all facts over  $sch(\mathcal{P})$  that are available at  $x_i$  during step  $s_i$ , i.e., the joint result of rules in  $pure(\mathcal{P})$  of the form (1), (2), and (6).

*Sending.* We define  $snd_{\mathcal{R}}^{[i]}$  to represent the sending of messages during transition  $i$ . We proceed as follows. Let  $mesg_{\mathcal{R}}^{[i]}$  denote the output of subprogram  $async_{\mathcal{P}}$  during transition  $i$ , restricted to the facts having their addressee-component in the network. Now, we define  $snd_{\mathcal{R}}^{[i]}$  to consist of the following facts:

- all facts  $cand_R(x_i, s_i, y, t, \bar{a})$  for which  $R(y, \bar{a}) \in mesg_{\mathcal{R}}^{[i]}$  and  $t \in \mathbb{N}$  such that  $(y, t) \not\prec_{\mathcal{R}}(x_i, s_i)$ , representing the result of rule (9);
- all facts  $chosen_R(x_i, s_i, y, t, \bar{a})$  for which  $R(y, \bar{a}) \in mesg_{\mathcal{R}}^{[i]}$  and  $t = loc_{\mathcal{R}}(j)$  with  $j = \alpha_{\mathcal{R}}(i, y, R(\bar{a}))$ , representing the result of rule (4); and,
- all facts  $other_R(x_i, s_i, y, u, \bar{a})$  for which  $R(y, \bar{a}) \in mesg_{\mathcal{R}}^{[i]}$ ,  $u \in \mathbb{N}$ ,  $(y, u) \not\prec_{\mathcal{R}}(x_i, s_i)$  and  $u \neq loc_{\mathcal{R}}(j)$  with  $j = \alpha_{\mathcal{R}}(i, y, R(\bar{a}))$ , representing the result of rule (5).

*Conclusion.* We can show that  $M$  is indeed a model of  $\mathcal{P}$  on input  $H$ ; this proof can be found in Appendix A of the online appendix to the paper. By construction of  $M$ , we have, as desired:

$$M|_{sch(\mathcal{P})^{LT}} = \bigcup_{i \in \mathbb{N}} duc_{\mathcal{R}}^{[i]} = \bigcup_{i \in \mathbb{N}} D_i^{\uparrow x_i, s_i} = trace(\mathcal{R}).$$

### 5.3 Model to run

Let  $\mathcal{P}$  be a Dedalus program and let  $H$  be an input distributed database instance for  $\mathcal{P}$ , over some network  $\mathcal{N}$ . Let  $M$  be a model of  $\mathcal{P}$  on input  $H$ . We show there is a fair run  $\mathcal{R}$  of  $\mathcal{P}$  on input  $H$  such that  $trace(\mathcal{R}) = M|_{sch(\mathcal{P})^{LT}}$ .

The direction shown in Section 5.2 is perhaps the most intuitive direction because we only have to show that a concrete set of facts is actually a stable model. In this section, we do not yet understand what  $M$  can contain. So, a first important step is to show that  $M$  has some desirable properties which allow us to construct a run from it.

Using the notation from Section 3.2.3, let  $G$  abbreviate the ground program  $ground_M(C, I)$  where  $C = pure(\mathcal{P})$  and  $I = decl(H)$ . By definition of  $M$  as a stable model, we have  $M = G(I)$ .

First, it is important to know that in  $M$  we find location specifiers where we expect location specifiers and we find timestamps where we expect timestamps. Formally, we call  $M$  *well-formed* if:

- for each  $R(x, s, \bar{a}) \in M|_{sch(\mathcal{P})^{LT}}$  we have  $x \in \mathcal{N}$  and  $s \in \mathbb{N}$ ;
- for each  $before(x, s, y, t) \in M$ , we have  $x, y \in \mathcal{N}$  and  $s, t \in \mathbb{N}$ ;
- for each fact  $cand_R(x, s, y, t, \bar{a})$ ,  $chosen_R(x, s, y, t, \bar{a})$  and  $other_R(x, s, y, t, \bar{a})$  in  $M$ , we have  $x, y \in \mathcal{N}$  and  $s, t \in \mathbb{N}$ ;
- for each fact  $hasSender(x, s, y, t)$ ,  $isSmaller(x, s, y, t)$ ,  $hasMax(x, s, y)$  and  $rcvInf(x, s)$  in  $M$ , we have  $x, y \in \mathcal{N}$  and  $s, t \in \mathbb{N}$ .

It can be shown by induction on the fixpoint computation of  $G$  that  $M$  is always well-formed. We omit the details.

The rest of this subsection is organized as follows. In Section 5.3.1, we extract a happens-before relation  $<_M$  from  $M$ . Next, in Section 5.3.2, we construct a run  $\mathcal{R}$ : we use  $<_M$  to establish a total order on  $\mathcal{N} \times \mathbb{N}$  that tells us which are the active nodes in the transitions of  $\mathcal{R}$ . Finally, we show in Section 5.3.3 that  $\mathcal{R}$  is fair.

### 5.3.1 Partial order

We define the following relation  $<_M$  on  $\mathcal{N} \times \mathbb{N}$ : for each  $(x, s) \in \mathcal{N} \times \mathbb{N}$  and  $(y, t) \in \mathcal{N} \times \mathbb{N}$ , we write  $(x, s) <_M (y, t)$  if and only if  $\text{before}(x, s, y, t) \in M$ . The rest of this section is dedicated to showing that  $<_M$  is a well-founded strict partial order on  $\mathcal{N} \times \mathbb{N}$ .

Let  $G$  abbreviate the ground program  $\text{ground}_M(C, I)$  where  $C = \text{pure}(\mathcal{P})$  and  $I = \text{decl}(H)$ . Regarding terminology, an edge  $(x, s) <_M (y, t)$  is called a *local edge*, a *message edge* or a *transitive edge* if the fact  $\text{before}(x, s, y, t) \in M$  can be derived by a ground rule in  $G$  of respectively the form (7), the form (10), or the form (8).<sup>15</sup> It is possible that an edge is of two or even three types at the same time.

Consider the following claim:

*Claim 5*

Relation  $<_M$  is a strict partial order on  $\mathcal{N} \times \mathbb{N}$ .

*Proof*

We show that  $<_M$  is transitive and irreflexive.

*Transitive.* First, we show that  $<_M$  is transitive. Suppose we have  $(x, s) <_M (z, u)$  and  $(z, u) <_M (y, t)$ . We have to show that  $(x, s) <_M (y, t)$ . By definition of  $<_M$ , we have  $\text{before}(x, s, z, u) \in M$  and  $\text{before}(z, u, y, t) \in M$ . Because rule (8) is positive, we have the following ground rule in  $G$ :

$$\text{before}(x, s, y, t) \leftarrow \text{before}(x, s, z, u), \text{before}(z, u, y, t).$$

Because  $M$  is a stable model and the body of the previous ground rule is in  $M$ , we obtain  $\text{before}(x, s, y, t) \in M$ . Hence,  $(x, s) <_M (y, t)$ , as desired.

*Irreflexive.* Because an edge  $(x, s) <_M (x, s)$  for any  $(x, s) \in \mathcal{N} \times \mathbb{N}$  would form a cycle of length one, it is sufficient to show that there are no cycles in  $<_M$  at all. This gives us irreflexivity, as desired.

First, let  $<'_M$  denote the restriction of  $<_M$  to the edges that are local or message edges. Note that this definition allows some edges in  $<'_M$  to also be transitive. The edges that are missing from  $<'_M$  with respect to  $<_M$  are only derivable by ground rules of the form (8); we call these the *pure* transitive edges. We start by showing that  $<'_M$  contains no cycles. We show this with a proof by contradiction. So, suppose that there is a cycle in  $\mathcal{N} \times \mathbb{N}$  through the edges of  $<'_M$ :

$$(x_1, s_1) <_M (x_2, s_2) <_M \dots <_M (x_n, s_n),$$

<sup>15</sup> The body of such a ground rule has to be in  $M$ .

with  $n \geq 2$  and  $(x_1, s_1) = (x_n, s_n)$ . We have  $\text{before}(x_i, s_i, x_{i+1}, s_{i+1}) \in M$  for each  $i \in \{1, \dots, n-1\}$ . Based on these *before*-facts, ground rules in  $G$  of the form (8) will have derived  $\text{before}(x_i, s_i, x_j, s_j) \in M$  for each  $i, j \in \{1, \dots, n\}$ .

If each edge on the above cycle would be only local, then for each  $i, j \in \{1, \dots, n\}$  with  $i < j$  we have  $x_i = x_j$  and  $s_i < s_j$ , and hence  $s_1 \neq s_n$ , which is false. So, there has to be some  $k \in \{1, \dots, n-1\}$  such that  $(x_k, s_k) <_M (x_{k+1}, s_{k+1})$  is a message edge, derived by a ground rule of the form (10):

$$\text{before}(x_k, s_k, x_{k+1}, s_{k+1}) \leftarrow \text{chosen}_R(x_k, s_k, x_{k+1}, s_{k+1}, \bar{a}).$$

Therefore,  $\text{chosen}_R(x_k, s_k, x_{k+1}, s_{k+1}, \bar{a}) \in M$ . This *chosen*<sub>R</sub>-fact must be derived by a ground rule of the form (4) in  $G$ , which implies that

$$\text{cand}_R(x_k, s_k, x_{k+1}, s_{k+1}, \bar{a}) \in M.$$

This *cand*<sub>R</sub>-fact must in turn be derived by a ground rule  $\psi$  of the form (9). Because rules of the form (9) in *pure*( $\mathcal{P}$ ) contain a negative *before*-atom in their body, the presence of  $\psi$  in  $G$  requires that  $\text{before}(x_{k+1}, s_{k+1}, x_k, s_k) \notin M$ . But that is a contradiction, because  $\text{before}(x_i, s_i, x_j, s_j) \in M$  for each  $i, j \in \{1, \dots, n\}$  (see above).

Now, we show there are no cycles in the entire relation  $<_M$ . Since  $M = G(\text{decl}(H))$ , we have  $M = \bigcup_{i \in \mathbb{N}} M_i$  where  $M_0 = \text{decl}(H)$  and  $M_i = T(M_{i-1})$  for each  $i \geq 1$  where  $T$  is the immediate consequence operator of  $G$ . By induction on  $i$ , we show that an edge  $\text{before}(x, s, y, t) \in M_i$  either is a local or message edge, or it can be replaced by a path of local or message edges in  $M_i$ . Then any cycle in  $<_M$  would imply there is a cycle in  $<'_M$ , which is impossible. So,  $<_M$  cannot contain cycles. Now, this induction property is satisfied for the base case because  $M_0$  does not contain *before*-facts. For the induction hypothesis, assume the property holds for  $M_{i-1}$ , where  $i \geq 1$ . For the inductive step, let  $\text{before}(x, s, y, t) \in M_i \setminus M_{i-1}$ . If this fact is derived by a ground rule of the form (7) or (10) then the property is satisfied. Now suppose the fact is derived by a ground rule of the form (8):

$$\text{before}(x, s, y, t) \leftarrow \text{before}(x, s, z, u), \text{before}(z, u, y, t).$$

Both body facts are in  $M_{i-1}$ , implying  $M_{i-1}$  contains a path of local or message edges from  $(x, s)$  to  $(z, u)$  and from  $(z, u)$  to  $(y, t)$ . Hence, using  $M_{i-1} \subseteq M_i$ , the edge  $\text{before}(x, s, y, t) \in M_i$  can be replaced by a path of local or message edges in  $M_i$ . □

In Section 4.5 we have added extra rules to *pure*( $\mathcal{P}$ ) to enforce that every node only receives a finite number of messages during each step. We now verify that this works correctly:

*Claim 6*

For each  $(y, t) \in \mathcal{N} \times \mathbb{N}$  there are only a finite number of pairs  $(x, s) \in \mathcal{N} \times \mathbb{N}$  such that  $(x, s) <_M (y, t)$  is a message edge.

*Proof*

We start by noting that  $M$  does not contain the fact  $\text{rcvInf}(y, t)$ . Indeed, in order to derive this fact, we need a ground rule in  $G$  of the form (14), which has a body fact

of the form  $\text{hasSender}(y, t, x, s)$ . Such  $\text{hasSender}$ -facts must be generated by ground rules in  $G$  of the form (11). The rule (11) negatively depends on relation  $\text{rcvInf}$ . Thus, specifically, if we want a ground rule in  $G$  that can derive  $\text{hasSender}(y, t, x, s)$ , we should require the absence of  $\text{rcvInf}(y, t)$  from  $M$ . So  $\text{rcvInf}(y, t) \in M$  requires  $\text{rcvInf}(y, t) \notin M$ , which is impossible.

The rest of the proof works towards a contradiction. So, suppose that  $(y, t)$  has an infinite number of incoming message edges. Because there are only a finite number of nodes in  $\mathcal{N}$ , there has to be a node  $x$  that has an infinite number of timestamps  $s$  such that  $\text{before}(x, s, y, t) \in M$  is a message edge. Since it is a message edge, such a fact  $\text{before}(x, s, y, t)$  can be generated by a ground rule in  $G$  of the form (10), which implies that there is a relation  $R$  in  $\text{idb}(\mathcal{P})$  and a tuple  $\bar{a}$  such that  $\text{chosen}_R(x, s, y, t, \bar{a}) \in M$ . Because  $\text{rcvInf}(y, t) \notin M$  (see above), for each of these  $\text{chosen}_R$ -facts, there is a ground rule of the form (11) in  $M$  that derives  $\text{hasSender}(y, t, x, s) \in M$ .

Rule (14) has a negative  $\text{hasMax}$ -atom in its body. If we can show that  $\text{hasMax}(y, t, x) \notin M$ , then there will be a ground rule in  $G$  of the form (14), where  $\text{hasSender}(y, t, x, s) \in M$ :

$$\text{rcvInf}(y, t) \leftarrow \text{hasSender}(y, t, x, s).$$

This then causes  $\text{rcvInf}(y, t) \in M$ , giving the desired contradiction.

Also towards a proof by contradiction, suppose that  $\text{hasMax}(y, t, x) \in M$ . This means that there is a ground rule  $\psi$  in  $G$  of the form (13):

$$\text{hasMax}(y, t, x) \leftarrow \text{hasSender}(y, t, x, s).$$

Because the rule (13) contains a negative  $\text{isSmaller}$ -atom in the body, and because  $\psi \in G$ , we know that  $\text{isSmaller}(y, t, x, s) \notin M$ . But because there are infinitely many facts of the form  $\text{hasSender}(y, t, x, s') \in M$ , there is at least one fact  $\text{hasSender}(y, t, x, s') \in M$  with  $s < s'$ . Moreover, the rule (12) is positive, and therefore the following ground rule is always in  $G$ :

$$\text{isSmaller}(y, t, x, s) \leftarrow \text{hasSender}(y, t, x, s), \text{hasSender}(y, t, x, s'), s < s'.$$

Since the body of this ground rule is in  $M$ , the rule derives  $\text{isSmaller}(y, t, x, s) \in M$ , which gives the desired contradiction. □

An ordering  $<$  on a set  $A$  is called *well-founded* if for each  $a \in A$ , there are only a finite number of elements  $b \in A$  such that  $b < a$ . We now use Claim 6 to show:

*Claim 7*

Relation  $<_M$  on  $\mathcal{N} \times \mathbb{N}$  is well-founded.

*Proof*

Let  $(x, s) \in \mathcal{N} \times \mathbb{N}$ . We have to show that there are only a finite number of pairs  $(y, t) \in \mathcal{N} \times \mathbb{N}$  such that  $(y, t) <_M (x, s)$ . Technically, we can limit our attention to paths in  $<_M$  consisting of local edges and message edges, because if we can show that there are only a finite number of predecessors of  $(x, s)$  on such paths, then there are only a finite number of predecessors when we include the transitive edges

as well. First, we show that every pair  $(y, t) \in \mathcal{N} \times \mathbb{N}$  has only a finite number of incoming local and message edges. If  $t > 0$ , we can immediately see that  $(y, t)$  has precisely one incoming local edge, as created by a ground rule of the form (7), and if  $t = 0$  then  $(y, t)$  has no incoming local edge. Also, Claim 6 tells us that  $(y, t)$  has only a finite number of incoming message edges. So, the number of incoming local and message edges in  $(y, t)$  is finite.

Let  $(y, t) \in \mathcal{N} \times \mathbb{N}$  be a pair such that  $(y, t) <_M (x, s)$  is a local edge or a message edge. Starting in  $(x, s)$ , we can follow this edge backwards so that we reach  $(y, t)$ . If  $(y, t)$  itself has incoming local or message edges, from  $(y, t)$  we can again follow an edge backwards. This way we can incrementally construct backward paths starting from  $(x, s)$ . Because at each pair of  $\mathcal{N} \times \mathbb{N}$  there are only a finite number of incoming local or message edges (shown above), if  $(x, s)$  would have an infinite number of predecessors, we must be able to construct a backward path of infinite length. We now show that the existence of such an infinite path leads to a contradiction. So, suppose that there is a backward path of infinite length. Because there are only a finite number of nodes in the network  $\mathcal{N}$ , there must be a node  $y$  that occurs infinitely often on this path. We will now show that, as we progress further along the backward path, we must see the local timestamps of  $y$  strictly decrease. Hence, we must eventually reach timestamp 0 of  $y$ , after which we cannot decrement the timestamps of  $y$  anymore, and thus it is impossible that  $y$  occurs infinitely often along the path. Suppose that the timestamps of  $y$  do not strictly decrease. There are two cases. First, if the same pair  $(y, t)$  would occur twice on the path, we would have a cycle in  $<_M$ , which is not possible by Claim 5. Secondly, suppose that there are two timestamps  $t$  and  $t'$  of  $y$  such that  $t < t'$  and  $(y, t)$  occurs before  $(y, t')$  on the backward path, meaning that  $(y, t)$  lies closer to  $(x, s)$ . Because the edges were followed in reverse, we have

$$(y, t') <_M \dots <_M (y, t).$$

But since  $t < t'$ , by means of local edges, we always have

$$(y, t) <_M (y, t + 1) <_M \dots <_M (y, t').$$

So, there would be a cycle between  $(y, t')$  and  $(y, t)$ . But that is again impossible by Claim 5. □

### 5.3.2 Construction of run

Let  $<_M$  be the well-founded strict partial order on  $\mathcal{N} \times \mathbb{N}$  as defined in the preceding subsection. The relation  $<_M$  has the intuition of a happens-before relation of a run (Section 5.2.1), but the novelty is that it comes from a purely declarative model  $M$ . We will now use  $<_M$  to construct a run  $\mathcal{R}$  such that  $trace(\mathcal{R}) = M|_{sch(\mathcal{P})^{LT}}$ .

*Total order.* It is well known that a well-founded strict partial order can be extended to a well-founded strict total order. So, let  $<_M$  be a well-founded strict total order on  $\mathcal{N} \times \mathbb{N}$  that extends  $<_M$ , i.e., for each  $(x, s) \in \mathcal{N} \times \mathbb{N}$  and  $(y, t) \in \mathcal{N} \times \mathbb{N}$ , if  $(x, s) <_M (y, t)$  then  $(x, s) <_M (y, t)$ , but the reverse does not have to hold.

Ordering the set  $\mathcal{N} \times \mathbb{N}$  according to  $<_M$  gives us a sequence of pairs that will form the transitions in the constructed run  $\mathcal{R}$ . Concretely, we obtain a sequence of nodes by taking the node-component from each pair. This will form our sequence of active nodes. Similarly, by taking the timestamp-component from each pair of  $\mathcal{N} \times \mathbb{N}$ , we obtain a sequence of timestamps. These are the local clocks of the active nodes during their transitions.

We introduce some extra notations to help us reason about the ordering of time that is implied by  $<_M$ . For each  $(x, s) \in \mathcal{N} \times \mathbb{N}$ , let  $glob_M(x, s) \in \mathbb{N}$  denote the ordinal of  $(x, s)$  as implied by  $<_M$ , which is well-defined because  $<_M$  is well-founded. For technical convenience, we let ordinals start at 0. Note,  $glob_M(\cdot)$  is an injective function. For any  $i \in \mathbb{N}$ , we define  $(x_i, s_i)$  to be the unique pair in  $\mathcal{N} \times \mathbb{N}$  such that  $glob_M(x_i, s_i) = i$ .

As a counterpart to function  $glob_M(\cdot)$ , for each  $i \in \mathbb{N}$  and each  $x \in \mathcal{N}$ , let  $loc_M(i, x)$  denote the *size* of the set

$$\{s \in \mathbb{N} \mid glob_M(x, s) < i\}.$$

Intuitively, if  $i$  is regarded to be the ordinal of a transition in a run,  $loc_M(i, x)$  is the number of local steps of  $x$  that came before transition  $i$ , i.e., the number of transitions before  $i$  in which  $x$  was the active node. If  $x = x_i$  (the active node) then  $loc_M(i, x)$  is effectively the timestamp of  $x$  *during* transition  $i$ , and if  $x \neq x_i$  then  $loc_M(i, x)$  is the next timestamp of  $x$  that still has to come *after* transition  $i$ . Note, the functions  $glob_M(\cdot)$  and  $loc_M(\cdot)$  closely resemble the functions  $glob_{\mathcal{R}}(\cdot)$  and  $loc_{\mathcal{R}}(\cdot)$  of Section 5.1.5.

*Configurations.* We will now define the desired run  $\mathcal{R}$  of  $\mathcal{P}$  on  $H$ . First, we define an infinite sequence of configurations  $\rho_0, \rho_1, \rho_2$ , etc. In a second step, we will connect each pair of subsequent configurations by a transition. Recall from Section 5.1.1 that a configuration describes for each node what facts it has stored locally (state), and also what messages have been sent to this node but that are not yet received (message buffer). The facts that are stored on a node are either input *edb*-facts, or facts derived by inductive rules in a previous step of the node. The first kind of facts can be easily obtained from  $M$  by keeping only the facts over schema  $edb(\mathcal{P})^{LT}$ , which gives a subset of  $decl(H)$ .

For the second kind of state facts, we look at the inductively derived facts in  $M$ . Rules in  $pure(\mathcal{P})$  that represent inductive rules of  $\mathcal{P}$  are recognizable as rules of the form (2): they have a head atom over  $sch(\mathcal{P})^{LT}$  and they have a (positive) *tsucc*-atom in their body. No other kind of rule in  $pure(\mathcal{P})$  has this form. Hence, the ground rules in  $G$  that are based on rules of the form (2) are also easily recognizable, and we will call these *inductive ground rules*. A ground rule  $\psi \in G$  is called *active* on  $M$  if  $pos_\psi \subseteq M$ , which implies  $head_\psi \in M$  because  $M$  is stable. Let  $M^{ind}$  denote all head atoms of inductive ground rules in  $G$  that are active on  $M$ . Note that  $M^{ind} \subseteq M$ . Regarding notation, for an instance  $I$  over  $sch(\mathcal{P})^{LT}$ , we write  $I^\Downarrow$  to denote the set  $\{R(\bar{a}) \mid \exists x, s : R(x, s, \bar{a}) \in I\}$ , and we write  $I|^{x,s}$  to denote the set  $\{R(y, t, \bar{a}) \in I \mid y = x, t = s\}$ .

Now, for each  $i \in \mathbb{N}$ , for each node  $x \in \mathcal{N}$ , denoting  $s = loc_M(i, x)$ , in configuration  $\rho_i = (st_i, bf_i)$ , the state  $st_i(x)$  is defined as

$$((M|_{edb(\mathcal{P})^{LT}})^{x,s} \cup M^{ind}|_{x,s})^\Downarrow.$$

We remove the location specifier and timestamp because we have to obtain facts over the schema of  $\mathcal{P}$ , not over the schema of  $pure(\mathcal{P})$ .

Now we define the message buffers in the configurations. Recall that the message buffer of a node always contains pairs of the form  $(j, \mathbf{f})$ , where  $j \in \mathbb{N}$  is the transition in which fact  $\mathbf{f}$  was sent. For each  $i \in \mathbb{N}$ , for each node  $x \in \mathcal{N}$ , in configuration  $\rho_i = (st_i, bf_i)$ , the message buffer  $bf_i(x)$  is defined as

$$\{(glob_M(y, t), R(\bar{a})) \mid \exists u : chosen_R(y, t, x, u, \bar{a}) \in M, glob_M(y, t) < i \leq glob_M(x, u)\}.$$

Note the use of addressee  $x$  in this definition. The definition of  $bf_i(x)$  reflects the operational semantics, in that the messages in the buffer of node  $x$  must be sent in a previous transition, as expressed by the constraint  $glob_M(y, t) < i$ . Moreover, the constraint  $i \leq glob_M(x, u)$  says that  $bf_i(x)$  contains only messages that will be delivered in transitions of  $x$  that come after configuration  $\rho_i$ . Possibly  $i = glob_M(x, u)$ , and in that case the message will be delivered in the transition immediately after configuration  $\rho_i$ , which is transition  $i$  (see also below).

*Transitions.* So far, we have obtained a sequence of configurations  $\rho_0, \rho_1, \rho_2$ , etc. Now we define a sequence of tuples, one tuple per ordinal  $i \in \mathbb{N}$ , that represents the transition  $i$ . Let  $i \in \mathbb{N}$ . Recall from above that  $(x_i, s_i)$  is the unique pair in  $\mathcal{N} \times \mathbb{N}$  such that  $glob_M(x_i, s_i) = i$ . The tuple  $\tau_i$  is defined as  $(\rho_i, x_i, m_i, i, \rho_{i+1})$ , where

$$m_i = \{(glob_M(y, t), R(\bar{a})) \mid chosen_R(y, t, z, u, \bar{a}) \in M, glob_M(z, u) = i\}.$$

Intuitively,  $m_i$  selects all messages that arrive in transition  $i$ . And since  $glob_M(z, u) = i$  implies  $z = x_i$  and  $u = s_i$ , we thus select all messages destined for step  $s_i$  of node  $x_i$ .

*Trace.* We can show that sequence  $\mathcal{R}$  is indeed a legal run of  $\mathcal{P}$  on input  $H$  such that  $trace(\mathcal{R}) = M|_{sch(\mathcal{P})^{LT}}$ ; this proof can be found in Appendix B of the online appendix to the paper. In the following subsection we show that  $\mathcal{R}$  is also fair.

### 5.3.3 Fair run

Let  $\mathcal{R}$  be the run as constructed in the previous subsection. We now show that  $\mathcal{R}$  is fair. For each transition index  $i \in \mathbb{N}$ , let  $\rho_i = (st_i, bf_i)$  denote the source configuration of transition  $i$ . Recall from Section 5.1.4 that we have to check two fairness conditions:

- (1) every node is the active node in an infinite number of transitions; and
- (2) for every transition  $i \in \mathbb{N}$ , for every  $y \in \mathcal{N}$ , for every pair  $(j, \mathbf{f}) \in bf_i(y)$ , there is a transition  $k$  with  $i \leq k$  in which  $(j, \mathbf{f})$  is delivered to  $y$ .

We show that  $\mathcal{R}$  satisfies the first fairness condition. Let  $x \in \mathcal{N}$  be a node, and let  $s \in \mathbb{N}$  be a timestamp of  $x$ . Consider transition  $i = glob_M(x, s)$ . This transition

has active node  $x_i = x$ . We can find such a transition with active node  $x$  for every timestamp  $s \in \mathbb{N}$  of  $x$ , and these transitions are all unique because function  $glob_M(\cdot)$  is injective. So, there are an infinite number of transitions in  $\mathcal{R}$  with active node  $x$ .

We show that  $\mathcal{R}$  satisfies the second fairness condition. Let  $i \in \mathbb{N}$ ,  $y \in \mathcal{N}$ , and  $(j, \mathbf{f}) \in bf_i(y)$ . Denote  $\mathbf{f} = R(\bar{a})$ . From its construction, the pair  $(j, \mathbf{f}) \in bf_i(y)$  implies there are values  $x \in \mathcal{N}$ ,  $s \in \mathbb{N}$  and  $t \in \mathbb{N}$  such that  $chosen_R(x, s, y, t, \bar{a}) \in M$  and  $j = glob_M(x, s) < i \leq glob_M(y, t)$ . Denote  $k = glob_M(y, t)$ . Hence,  $i \leq k$  and  $(j, \mathbf{f}) \in m_k$  by definition of  $m_k$ . Thus,  $(j, \mathbf{f})$  is delivered to  $x_k = y$  in transition  $k$ .

## 6 Discussion

We have represented distributed programs in Datalog under the stable model semantics. Moreover, we have shown that the stable models represent the desired behavior of the distributed program, as found in a realistic operational semantics. We now discuss some points for future work.

As mentioned, many Datalog-inspired languages have been proposed to implement distributed applications (Loo et al. 2009; Navarro and Rybalchenko 2009; Grumbach and Wang 2010; Abiteboul et al. 2011), and they contain several powerful features such as aggregation and non-determinism (choice). Our current framework already represents the essential features that all these languages possess: reasoning about distributed state and representing message sending. Nonetheless, we have probably not yet explored the full power of stable models. We therefore expect that this work can be extended to languages that incorporate more powerful language constructs such as the ones mentioned above. It might also be possible to remove the syntactic stratification condition that we have used for the deductive rules.

More related to multi-agent systems (Leite et al. 2002; Nigam and Leite 2006; Leite and Soares 2007), it might be interesting to allow logic programs used in declarative networking to dynamically modify their rules. The question would be how (and if) this can be represented in our model-based semantics.

The effect of variants of the model-based semantics can be studied. For example, messages can be sent into the past when the causality rules are removed. Then, one might ask which (classes of) programs still work “correctly” under such a non-causal semantics; some preliminary results are in Ameloot and den Bussche (2014).

Lastly, we can think about the output of distributed Datalog programs. Marczak et al. (2011) define the output with *ultimate* facts, which are facts that will eventually always be present on the network. This way, the output of a run (or equivalently stable model) can be defined. Then, a *consistent* program is required to produce the same output in every run. For consistent programs, the output on an input distributed database instance can thus be defined as the output of any run. We can now consider the following decision problem: for a consistent program, an input distributed database instance for that program, and a fact, decide if this fact is output by the program on that input. We think that decidability depends on the semantics of the message buffers. In this paper, we have represented per addressee duplicate messages in its message buffer. This is a realistic representation, since in a real network, the same message can be sent multiple times, and hence,

multiple instances of the same message can be in transmission simultaneously. If we would forbid duplicate messages in the buffers, then the decision problem becomes decidable because only a finite number of configurations would be possible by finiteness of the input domain. But when duplicates are preserved, the number of configurations is not limited, and we expect that the problem will be undecidable in general. However, we might want to investigate whether decidability can be obtained in particular (syntactically defined) cases. If so, it might be interesting for those cases to find finite representations of the stable models. This could serve as a more intuitive programmer abstraction, or it could perhaps be used to more efficiently simulate the behavior of the network for testing purposes.

### Acknowledgment

The second author thanks Serge Abiteboul for a number of interesting discussions.

### Supplementary material

For supplementary material for this article, please visit <http://dx.doi.org/10.1017/S1471068415000381>

### References

- ABITEBOUL, S., BIENVENU, M., GALLAND, A., *et al.* 2011. A rule-based language for Web data management. In *Proc. 30th ACM Symposium on Principles of Database Systems*, ACM Press, New York, NY, USA, 293–304.
- ABITEBOUL, S., HULL, R. AND VIANU, V. 1995. *Foundations of Databases*. Addison-Wesley.
- ABRIAL, J. 2010. *Modeling in Event-B – System and Software Engineering*. Cambridge University Press.
- ALFERES, J., PEREIRA, L., PRZYMUSINSKA, H. AND PRZYMUSINSKI, T. 2002. LUPS—a language for updating logic programs. *Artificial Intelligence* 138, 1–2, 87–116.
- ALVARO, P., CONDIE, T., CONWAY, N., HELLERSTEIN, J. AND SEARS, R. 2009. I do declare: Consensus in a logic language. *Operating Systems Review* 43, 4, 25–30.
- ALVARO, P., CONWAY, N., HELLERSTEIN, J. AND MARCZAK, W. 2011. Consistency analysis in Bloom: A CALM and collected approach. In *Proc. 5th Biennial Conference on Innovative Data Systems Research*, 249–260. URL: [www.cidrdb.org](http://www.cidrdb.org).
- ALVARO, P., MARCZAK, W., *et al.* 2009. *Dedalus: Datalog in Time and Space*. Technical Report EECS-2009-173, University of California, Berkeley.
- ALVARO, P., MARCZAK, W. *et al.* 2011. Dedalus: Datalog in time and space. See de Moor *et al.* (2011), 262–281.
- AMELOOT, T. AND VAN DEN BUSSCHE, J. 2014. Positive Dedalus programs tolerate non-causality. *Journal of Computer and System Sciences* 80, 7, 1191–1213.
- AMELOOT, T., NEVEN, F. AND VAN DEN BUSSCHE, J. 2011. Relational transducers for declarative networking. In *Proc. 30th ACM Symposium on Principles of Database Systems*, ACM Press, 283–292.
- AMELOOT, T., NEVEN, F. AND VAN DEN BUSSCHE, J. 2013. Relational transducers for declarative networking. *Journal of the ACM* 60, 2, 15:1–15:38.

- AMELOOT, T. AND VAN DEN BUSSCHE, J. 2012. Deciding eventual consistency for a simple class of relational transducer networks. In *Proc. of the 15th International Conference on Database Theory*, ACM Press, 86–98.
- APT, K. AND BOL, R. 1994. Logic programming and negation: A survey. *The Journal of Logic Programming* 19-20, Supplement 1, 0, 9–71.
- APT, K., FRANCEZ, N. AND KATZ, S. 1988. Appraising fairness in languages for distributed programming. *Distributed Computing* 2, 226–241.
- ATTIYA, H. AND WELCH, J. 2004. *Distributed Computing: Fundamentals, Simulations, and Advanced Topics*. Wiley.
- CAVAGE, M. 2013. There's just no getting around it: You're building a distributed system. *ACM Queue* 11, 4.
- DE MOOR, O., GOTTLÖB, G., FURCHE, T. AND SELLERS, A. (Eds.) 2011. *Datalog Reloaded: First International Workshop, Datalog 2010*, Lecture Notes in Computer Science, Springer-Verlag, Berlin/Heidelberg, Germany, vol. 6702.
- DEUTSCH, A., SUI, L., VIANU, V. AND ZHOU, D. 2006. Verification of communicating data-driven Web services. In *Proc. 25th ACM Symposium on Principles of Database Systems*, ACM Press, 90–99.
- FRANCEZ, N. 1986. *Fairness*. Springer-Verlag New York, Inc., New York, NY, USA.
- GELFOND, M. AND LIFSCHITZ, V. 1988. The stable model semantics for logic programming. In *Proc. of the 5th International Conference on Logic Programming*, MIT Press, 1070–1080.
- GRUMBACH, S. AND WANG, F. 2010. Netlog, a rule-based language for distributed programming. In *Proc. 12th International Symposium on Practical Aspects of Declarative Languages*, M. Carro and R. Peña, Eds. Lecture Notes in Computer Science, vol. 5937, 88–103.
- HELLERSTEIN, J. 2010a. Datalog redux: Experience and conjecture. Video available (under the title “The Declarative Imperative”) from URL: <http://db.cs.berkeley.edu/jmh/>. PODS 2010 keynote.
- HELLERSTEIN, J. 2010b. The declarative imperative: Experiences and conjectures in distributed logic. *SIGMOD Record* 39, 1, 5–19.
- HUANG, S., GREEN, T. AND LOO, B. 2011. Datalog and emerging applications: An interactive tutorial. In *Proc. of the 2011 ACM SIGMOD International Conference on the Management of Data*, SIGMOD '11, ACM, 1213–1216.
- INTERLANDI, M., TANCA, L. AND BERGAMASCHI, S. 2013. Datalog in time and space, synchronously. In *Proc. 7th Alberto Mendelzon International Workshop on Foundations of Data Management*.
- JIM, T. 2001. SD3: A trust management system with certified evaluation. In *Proc. of the 2001 IEEE Symposium on Security and Privacy*, SP, IEEE Computer Society, 106–115.
- KRISHNAMURTHY, R. AND NAQVI, S. 1988. Non-deterministic choice in Datalog. In *Proc. of the 3rd International Conference on Data and Knowledge Bases*, Morgan Kaufmann, Burlington, MA, USA, 416–424.
- LAMPOR, L. 2000a. Distributed algorithms in TLA (abstract). In *Proc. of the 19th ACM Symposium on Principles of Distributed Computing*, ACM Press, 3.
- LAMPOR, L. 2000b. Fairness and hyperfairness. *Distributed Computing* 13, 4, 239–245.
- LEITE, J., ALFERES, J. AND PEREIRA, L. 2002. Minerva – a dynamic logic programming agent architecture. In *Revised Papers from the 8th International Workshop on Intelligent Agents VIII*, ATAL, Springer-Verlag, 141–157.
- LEITE, J. AND SOARES, L. 2007. Adding evolving abilities to a multi-agent system. In *Proc. of the 7th International Conference on Computational Logic in Multi-agent Systems*, CLIMA VII'06, Springer-Verlag, 246–265.

- LOBO, J., MA, J., RUSSO, A. AND LE, F. 2012. Declarative distributed computing. In *Correct Reasoning - Essays on Logic-Based AI in Honour of Vladimir Lifschitz*, E. Erdem, J. Lee, Y. Lierler and D. PEARCE, Eds., Lecture Notes in Computer Science, vol. 7265. Springer, 454–470.
- LOO, B. *et al.* 2009. Declarative networking. *Communications of the ACM* 52, 11, 87–95.
- LYNCH, N. 1996. *Distributed Algorithms*. Morgan Kaufmann.
- MA, J., LE, F., WOOD, D., RUSSO, A. AND LOBO, J. 2013. A declarative approach to distributed computing: Specification, execution and analysis. *Theory and Practice of Logic Programming* 13, Special Issue 4–5, 815–830.
- MARCZAK, W., ALVARO, P., CONWAY, N., HELLERSTEIN, J. AND MAIER, D. 2011. Confluence analysis for distributed programs: A model-theoretic approach. Technical Report UCB/EECS-2011-154 (Dec), EECS Department, University of California, Berkeley.
- MAREK, V. AND TRUSZCZYNSKI, M. 1999. Stable models and an alternative logic programming paradigm. In *The Logic Programming Paradigm*, K. Apt, V. Marek, M. Truszczynski and D. Warren, Eds., Artificial Intelligence, Springer, Berlin Heidelberg, 375–398.
- NAVARRO, J. AND RYBALCHENKO, A. 2009. Operational semantics for declarative networking. In *Proceedings 11th International Symposium on Practical Aspects of Declarative Languages*, A. Gill and T. Swift, Eds. Lecture Notes in Computer Science, 76–90.
- NIGAM, V. AND LEITE, J. 2006. A dynamic logic programming based system for agents with declarative goals. In *Proc. of the 4th International Conference on Declarative Agent Languages and Technologies*, DALT, Springer-Verlag, 174–190.
- SACCÀ, D. AND ZANIOLO, C. 1990. Stable models and non-determinism in logic programs with negation. In *Proc. of the 9th ACM Symposium on Principles of Database Systems*, ACM Press, 205–217.
- VARDI, M. 1982. The complexity of relational query languages. In *Proc. 14th ACM Symposium on the Theory of Computing*, ACM Press, New York, NY, USA, 137–146.
- WOODCOCK, J. AND DAVIES, J. 1996. *Using Z: Specification, Refinement, and Proof*. Prentice Hall.
- ZHANG, Q., CHENG, L. AND BOUTABA, R. 2010. Cloud computing: State-of-the-art and research challenges. *Journal of Internet Services and Applications* 1, 1, 7–18.
- ZINN, D., GREEN, T. AND LUDAESCHER, B. 2012. Win-move is coordination-free (sometimes). In *Proc. of the 15th International Conference on Database Theory*, ACM Press, 99–113.