

Санкт-Петербургский государственный университет  
Математическое обеспечение и администрирование информационных  
систем

Кафедра системного программирования

Намаконов Егор Сергеевич

# Компиляция модели памяти OCaml в Power

Магистерская диссертация

Научный руководитель:  
д.ф.-м.н. Кознов Д. В.

Консультант:  
к.ф.-м.н. Подкопаев А.В.

Рецензент:  
к.ф.-м.н. Березун Д. А.

Санкт-Петербург  
2020

ST. PETERSBURG UNIVERSITY  
Software and Administration of Information Systems

Software Engineering

Egor Namakonov

# Compilation of OCaml memory model to Power

Graduation Thesis

Scientific supervisor:  
D. Sc. Dmitri Koznov

Consultant:  
Ph.D. Anton Podkopaev

Reviewer:  
Ph.D. Daniil Berezun

St. Petersburg  
2020

# Оглавление

<b>1. Введение</b>	<b>4</b>
<b>2. Постановка задачи</b>	<b>6</b>
<b>3. Обзор</b>	<b>7</b>
3.1. Пример исполнения в слабой модели памяти . . . . .	7
3.2. Проблема корректности компиляции на примерах . . .	8
3.3. Графы исполнения . . . . .	12
3.4. Описание используемых моделей памяти . . . . .	15
3.4.1. Модель памяти OCaml (OCamlMM) . . . . .	15
3.4.2. Промежуточная модель памяти . . . . .	15
<b>4. Обоснование схемы компиляции</b>	<b>18</b>
4.1. Необходимость введения дополнительных инструкций	18
4.2. Схема компиляции . . . . .	20
<b>5. Доказательство корректности компиляции</b>	<b>21</b>
5.1. Соответствие графов исполнения . . . . .	21
5.2. Построение графа, соответствующего данному . . . . .	22
5.3. Доказательство последовательной согласованности по от- дельным адресам . . . . .	24
5.4. Доказательство отсутствия буферизации при чтении . .	25
<b>6. Формализация доказательства в Coq</b>	<b>30</b>
<b>7. Существующие исследования</b>	<b>33</b>
<b>8. Заключение</b>	<b>35</b>
<b>Список литературы</b>	<b>36</b>

# 1. Введение

Результат исполнения многопоточной программы, как правило, является недетерминированным. Конкретное множество допустимых результатов многопоточной программы определяется *моделью памяти* языка программирования. Наиболее известной является *модель последовательной согласованности* (sequential consistency, SC [11]). Она предполагает, что любой результат исполнения программы может быть получен путём попеременного исполнения инструкций отдельных потоков согласно программному порядку в них. Однако из-за оптимизаций, выполняемых современными компиляторами и процессорами, могут наблюдаться сценарии поведения, невозможные в такой модели. Так, на архитектуре x86 чтение по адресу в памяти может вернуть не самое последнее записанное значение, так как операция записи может быть буферизована.

Отказ от подобных оптимизаций нежелателен, поэтому современные модели памяти допускают некоторые сценарии поведения, невозможные в модели SC. Такие модели памяти называются *слабыми*. Например, слабыми являются модели памяти языков C++ [2], JavaScript [22] и Java [13], а также архитектур Power [1], x86 [16] и ARM [19].

Модель памяти OCam1 [4] (далее — OCam1MM) отличается свойством т.н. *локальной свободы от гонок по данным* (local data race freedom). Именно, гарантируется, что выполнение конфликтующих обращений по выбранному адресу в памяти (т.е. ситуация гонки по данным) не влияет на обращения к другим адресам, а также на последующие обращения по тому же адресу. Благодаря этому даже при возникновении гонки по данным в некоторый момент исполнения следующие участки программы будут исполнены согласно модели SC.

Чтобы гарантировать выполнение этого свойства, при компиляции нужно запретить некоторые оптимизации в зависимости от целевой архитектуры. Для этого может понадобиться, например, добавить в ассемблерный код инструкции-*барьеры*, которые запрещают нежелательные оптимизации на уровне процессора. Набор таких правил, по-

крывающий все возможные типы инструкций, называется *схемой компиляции*. Схема компиляции должна быть *корректной* — при исполнении любой программы, полученной при компиляции согласно этой схеме, должно наблюдаться только сценарии поведения, разрешённые OCamlMM для исходной программы.

Авторы OCamlMM разработали схемы компиляции OCamlMM в модели x86-TSO и ARMv8 [4] и доказали их корректность. При этом отсутствует схема компиляции в модель архитектуры Power. А между тем данная архитектура часто используется в современном серверном оборудовании [7]. Задача построения такой схемы осложнена тем, что модель Power, в отличие от моделей x86-TSO, ARMv8 и OCamlMM, не обладает т.н. свойством *multicopy atomicity*. Такое свойство означает, что записанные в память значения становятся доступны всем потокам в одном и том же порядке [19]. Из-за отсутствия этого свойства корректная схема компиляции OCamlMM в Power должна расставлять барьеры в результирующей программе так, чтобы запретить нежелательные сценарии поведения.

В рамках данной работы была поставлена задача разработать схему компиляции OCamlMM в модель Power и доказать её корректность. Для этого было решено использовать промежуточную модель памяти (Intermediate Memory Model, далее — IMM) [18], для которой уже доказана корректность компиляции в модель Power. Использование IMM как промежуточного этапа компиляции позволяет разбить доказательство корректности на два, которые впоследствии можно использовать в других доказательствах. Таким образом, построение схемы компиляции OCamlMM в IMM даёт схемы компиляции OCamlMM не только в Power, но и другие архитектуры, в которые компилируется IMM (на данный момент — x86 и ARM).

## 2. Постановка задачи

Целью данной работы является доказательство корректности компиляции модели памяти OCaml (OCamlMM) в модель памяти Power [1].

В работе были поставлены следующие задачи:

- построение схемы компиляции OCamlMM в IMM (для которой корректность компиляции в Power уже доказана);
- доказательство корректности полученной схемы;
- формализация доказательства в системе интерактивного доказательства теорем Coq [20].

### 3. Обзор

В этом разделе приводится пример слабого поведения программы и объясняются его причины. Затем на примерах рассматривается понятие корректности компиляции для моделей памяти. Далее формально описывается декларативный способ задания модели памяти [1], основанный на понятии графов исполнения. Наконец, формально описываются модели памяти OSamLMM и IMM, используемые далее в работе.

#### 3.1. Пример исполнения в слабой модели памяти

Рассмотрим программу, представленную на рис. 1. Здесь и далее используется упрощённый синтаксис программ:  $x$  и  $y$  обозначают адреса в памяти,  $a$  и  $b$  — локальные переменные (регистры),  $r\ l\ x$  — режим доступа (это понятие будет рассмотрено ниже). Сверху указаны начальные значения в памяти. В комментариях указаны наблюдаемые при чтении значения. Согласно модели SC, в зависимости от порядка исполнения инструкций в  $a$  и  $b$  могут быть записаны значения  $(1, 1)$ ,  $(1, 0)$  или  $(0, 1)$ . Однако после компиляции C++-аналога этой программы с помощью компилятора gcc и исполнения на архитектуре x86 в переменные  $a$  и  $b$  могут быть записаны нули, что не допускается моделью SC. У такого сценария поведения могут быть две причины. Во-первых, gcc может поменять местами обращения по разным адресам во время компиляции. Во-вторых, при исполнении на x86 возможна буферизация записи: в целях оптимизации обращений к памяти запись может быть отложена.

$$\frac{x = 0, y = 0}{\begin{array}{l|l} [x]^{r\ l\ x} := 1; & [y]^{r\ l\ x} := 1; \\ a := [y]^{r\ l\ x}; //0 & b := [x]^{r\ l\ x}; //0 \end{array}}$$

Рис. 1: Пример программы и её исполнения при буферизации записи

Новые сценарии поведения программы, возникающие в результате оптимизаций, могут быть некорректными с точки зрения требова-

ний к программе. Поэтому слабая модель памяти должна предоставлять возможность отменить оптимизации для отдельных инструкций. Для этого используются *режимы доступа*, различные по степени строгости. При компиляции инструкций с более строгими режимами компилятор может отменить некоторые оптимизации, а также добавить в результирующую программу инструкции-барьеры.

Так, в программе на рис. 1 был использован режим доступа `rlx`, который не ограничивает оптимизации соответствующих инструкций. Модель памяти C++ гарантирует, что если в этой программе для всех инструкций установить режим доступа `sc` вместо `rlx`, то поведение полученной программы будет согласовано с моделью SC. Это справедливо в силу того, что такие обращения будут скомпилированы с использованием инструкции MFENCE [14] — барьера памяти, запрещающего перестановки SC инструкций.

## 3.2. Проблема корректности компиляции на примерах

Модели OSam<sub>IMM</sub> и IMM определены декларативно [1]. Это означает, что каждое возможное поведение программы задаётся в виде *графа исполнения*, а семантика программы определяется как множество графов, удовлетворяющих некоторому условию. Пример программы и одного из графов её исполнения приведён на рис. 2.

Вершины графа соответствуют событиям — операциям над разделяемой памятью, которые производятся при выполнении инструкций программы. Так, событие  $w^{at}(x, 1)$  соответствует записи по адресу  $x$  значения 1 в режиме `at`; другими типами событий являются чтение и барьер памяти (обозначаются R и F соответственно). Кроме того, в графе выделяются инициализирующие события, которые соответствуют инициализирующей записи нулей в память. На рис. 2 все они для краткости обозначены множеством `Init`; далее в графах мы будем опускать эти события, если это не будет важно для рассуждений. Заметим, что содержимое локальных переменных потока не отражается в гра-

фе в явном виде, т.к. взаимодействие между потоками производится только через разделяемую память.

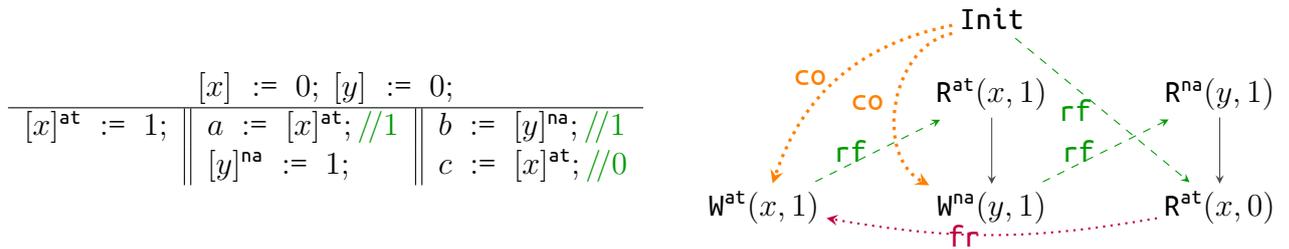


Рис. 2: Пример программы и сценария её исполнения, не согласованного в OCamlMM

Рёбра графа задают бинарные отношения между событиями. В данном графе есть четыре различных отношения: рёбра  $ro$  соответствуют программному порядку инструкций,  $rf$  — чтению записанного ранее значения,  $co$  — порядку выполнения записи по одному адресу,  $fr$  — чтению до указанного события записи. Отношения  $ro$  и  $co$  являются транзитивными, поэтому для их задания достаточно указывать только непосредственные рёбра. Кроме того, для краткости будем опускать подпись  $ro$  рядом с соответствующими рёбрами.

*Согласованными* (допустимыми моделью) называются те сценарии исполнения программы, графы которых удовлетворяют некоторому предикату, заданному моделью. В частности, предикат согласованности OCamlMM требует, чтобы в графе не было циклов, состоящих только из рёбер  $co$  и  $fr$ , проходящих между вершинами с меткой  $at$ , а также рёбер  $ro$  и  $rf$ . Это условие формализует свойство multicopy atomicity, описанное выше.

Граф исполнения на рис. 2 не является согласованным по OCamlMM. Действительно, этот сценарий исполнения нарушает свойство multicopy atomicity: второй поток читает записанное в  $x$  значение 1 до записи 1 в  $y$ , однако третий поток читает старое значение 0 из  $x$  после чтения 1 из  $y$ . Соответствующий граф исполнения не удовлетворяет предикату согласованности OCamlMM, так как между вершинами есть цикл, подходящий под описание выше. Таким образом, в OCamlMM после исполнения программы на рис. 2 переменные  $a$ ,  $b$  и  $c$  не могут содержать значения 1, 1 и 0 соответственно.

Условие корректности компиляции требует, чтобы сценарии поведения, запрещённые для исходной программы в ОСамЛММ, также были запрещены для скомпилированной программы в ИММ. Для декларативных моделей памяти это означает, что из несогласованности графа исполнения в ОСамЛММ должна следовать несогласованность соответствующего ему графа исполнения в ИММ. При этом, как будет показано далее, можно рассматривать вопрос согласованности по ОСамЛММ только для графа исполнения скомпилированной программы.

На рис. 3 приведён результат компиляции программы на рис. 2 согласно тривиальной схеме компиляции. Такая схема лишь заменяет режимы инструкций на их аналоги в ИММ:  $na$  заменяется на  $r\!l\!x$ , а  $at$  — на  $sc$ ; дополнительных инструкций не вводится. Соответственно, граф исполнения на рис. 3 отличается от графа на рис. 2 только метками вершин, и в нём сохраняется цикл того же вида. ИММ не гарантирует свойство *multicore atomicity*, и потому предикат её согласованности не требует отсутствия таких циклов. Поэтому граф на рис. 3 согласован, что делает соответствующий ему сценарий поведения разрешаемым ИММ, в отличие от ОСамЛММ. Поэтому тривиальная схема компиляции не является корректной.

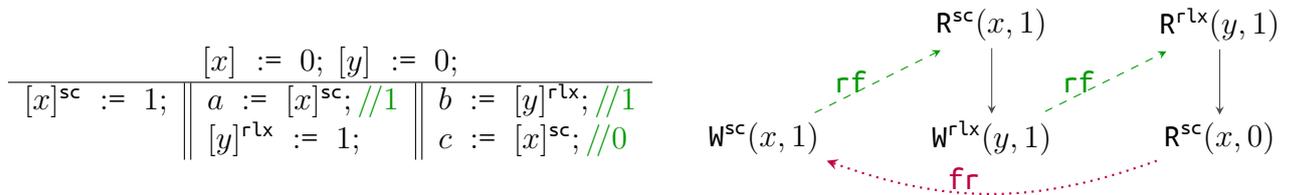


Рис. 3: Результат компиляции программы на рис. 2 с использованием тривиальной схемы компиляции и согласованный по ИММ граф его исполнения

На рис. 4 приведена программа, полученная в результате компиляции программы на рис. 2 согласно схеме компиляции, приведённой в разделе 4. В результате компиляции в ней появляются инструкции барьеров памяти, запрещающие некоторые оптимизации процессора и компилятора. С этими барьерами граф исполнения перестаёт быть согласованным: из рёбер  $rf$  между событиями с меткой  $sc$ , а также  $fr$ , ро и окружённого барьерами  $rf$  образуется цикл, запрещённый в ИММ.

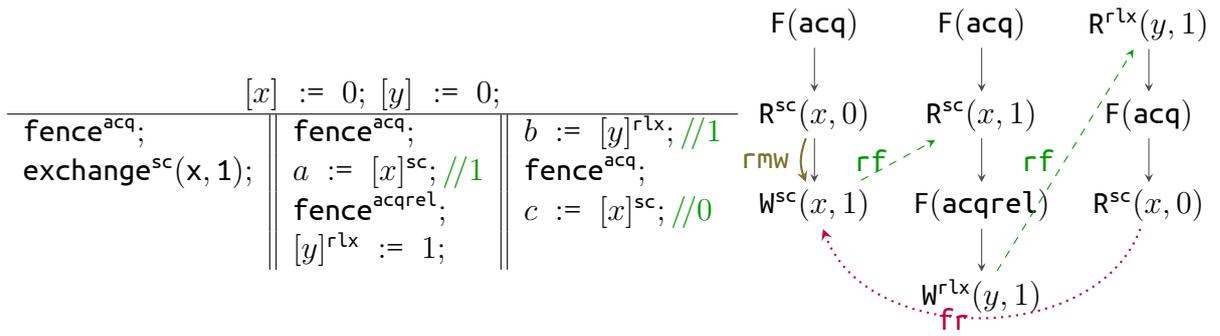


Рис. 4: Результат компиляции программы на рис. 2 с использованием схемы компиляции из раздела 4 и его граф исполнения, не согласованный в IMM

Таким образом, для доказательства корректности компиляции необходимо доказать следующую теорему<sup>1</sup>.

**Теорема 3.1.** Пусть  $PO = \parallel_{\tau \in \text{Tid}} PO_{\tau}$  и  $PI = \parallel_{\tau \in \text{Tid}} PI_{\tau}$  — программы для моделей OSamLM и IMM соответственно, причём для любого  $\tau \in \text{Tid}$  подпрограмма  $PI_{\tau}$  получена компиляцией  $PO_{\tau}$  с помощью схемы компиляции из таб. 1. Пусть  $G_I$  — согласованный по IMM граф исполнения  $PI$ . Тогда существует  $G_O$  — согласованный по OSamLM граф исполнения  $PO$ , соответствующий  $G_I$ .

Ключевой идеей доказательства является то, что согласованность по OSamLM можно рассматривать для обоих графов исполнения. Для доказательства теоремы достаточно показать, что  $G_I$  является согласованным по OSamLM. Из этого следует согласованность  $G_O$  по OSamLM, так как он фактически является подграфом  $G_I$ , а условия согласованности по OSamLM таковы, что выполняются для подграфов. Условие согласованности по OSamLM состоит в иррефлексивности одного отношения и ацикличности другого. Для каждого из этих отношений доказывается включение в такое отношение, для которого соответствующее условие выполняется в согласованном по IMM графе.

<sup>1</sup>Формальное понятие соответствие графов вводится в разделе 4.

### 3.3. Графы исполнения

В описаниях декларативных моделей памяти мы будем использовать следующие обозначения отношений между вершинами. Для бинарного отношения  $R$  обозначения  $R^?$ ,  $R^+$ ,  $R^*$  соответствуют его рефлексивному, транзитивному и транзитивно-рефлексивному замыканиям соответственно. Обратное отношение записывается как  $R^{-1}$ . Левая композиция отношений  $R_1$  и  $R_2$  записывается следующим образом:

$$R_1; R_2 \triangleq \{x, y | \exists z. (x, z) \in R_1 \wedge (z, y) \in R_2\}.$$

Непосредственные рёбра  $R$  обозначаются как  $R|_{\text{imm}} \triangleq R \setminus R; R$ . Тожественное отношение на множестве  $A$  обозначается как  $[A]$ ; в частности,  $[A]; R; [B] = R \cap (A \times B)$ .

В данном разделе описываются графы исполнения наиболее общего вида, без привязки к конкретным моделям памяти или языкам.

Считаем, что анализируемая программа  $P$  состоит из последовательных подпрограмм отдельных потоков  $P_\tau$ :  $P = \parallel_{\tau \in \text{Tid}} P_\tau$ , где  $\parallel$  — оператор параллельной композиции программ, а  $\text{Tid}$  — конечное множество идентификаторов потоков.

**Определение 3.2.** Граф исполнения  $G$  задаётся множеством вершин, бинарными отношениями на вершинах, а также функцией, сопоставляющей вершинам метки.

Множество вершин, обозначаемое как  $G.E$ , делится на инициализирующие события вида  $\text{Init } loc$  и неинициализирующие события вида  $\text{ThreadEvent } \tau n$ . Их компонентами являются:

- $loc \in \text{Loc}$  — адрес инициализации, где  $\text{Loc}$  — конечное множество адресов;
- $\tau \in \text{Tid}$  — номер потока;
- $n \in \mathbb{N}$  — порядковый номер внутри потока.

Функция  $G.\text{Lab}$  сопоставляет событиям метки вида  $(type, loc, mode, val)$ . Их компонентами являются:

- $type \in \{R, W, F\}$  — тип операции (чтение, запись, барьер);
- $loc \in Loc$  — адрес памяти (для барьера не определено);
- $mode$  — один из режимов доступа (например,  $rel$ ), частично упорядоченных отношением “строже чем” ( $\sqsubset$ ); конкретное множество режимов и их порядок определяется моделью памяти;
- $val \in Val$  — прочитанное/записанное значение (в случае барьера не определено), где  $Val$  — множество значений, которые могут храниться в памяти.

Следует отметить, что инициализирующие события вида  $Init\ loc$  обрабатываются особым образом. Именно,

$G.Lab(Init\ loc) = (W, loc, mode_{init}, val_{init})$ , где  $mode_{init}$  — выбранный режим доступа для инициализирующих событий записей (например, в IMM —  $rlx$ ), а  $val_{init}$  — начальное значение в памяти (как правило, 0).

Для множеств событий с определёнными метками вводятся соответствующие обозначения. Например, события с меткой чтения в режиме  $acq$  или более строгим будем обозначать как  $G.R^{acq}$  (или просто  $R^{acq}$ , если граф очевиден из контекста).

Рёбра графа представляют собой следующие отношения между событиями:

- программный порядок (program order):  $G.po(x, y) \iff (x \in Init \wedge y \notin Init) \vee (x.\tau = y.\tau \wedge x.n < y.n)$ ;
- порядок согласованности (coherence order):  $G.co = \bigcup_{l \in Loc} co_l$ , где  $co_l$  — тотальный порядок на событиях записи по адресу  $l$ ;
- наблюдение записанного значения (“читает-из”, reads from):

$$G.rf \subseteq \bigcup_{l \in Loc} G.W_l \times G.R_l, \text{ где}$$

$$G.rf(w, r) \Rightarrow G.Lab(w).val = G.Lab(r).val, \text{ codom}(G.rf) = G.R \text{ и}$$

$G.rf^{-1}$  является функциональным отношением;

- чтение до указанной записи:  $G.\mathbf{fr} = G.\mathbf{rf}^{-1}$ ;  $G.\mathbf{co}$  (from-read, “читает-до”).

Различные модели памяти могут иметь в графе исполнения и другие отношения. Например, в модели IMM также есть отношение  $\mathbf{rww} \subseteq \bigcup_{l \in \text{Loc}} [G.\mathbf{R}_l]; \text{po}|_{\text{imm}}; [G.\mathbf{W}_l]$ , соответствующее паре событий чтения и записи в операции read-modify-write.

Введём понятия сужения графа на поток  $i$ :  $G_\tau.E = \{e \in G.E \mid e.\tau = i\}$ ,  $G_\tau.\text{Lab} = G.\text{Lab}$ .

**Определение 3.3.** Графом исполнения программы  $P$  называется такой граф  $G$ , что его сужение на любой поток  $\tau \in \text{Pid}$  является однопоточным графом исполнения программы  $P_\tau$ .

Соответствие подпрограммы потока и однопоточного графа исполнения определяется средствами операционной семантики, специфичной для модели памяти [4], [18]. Мы не приводим подробностей здесь, скажем лишь, что такая семантика задаёт соответствие между выполнением инструкций языка и изменением графа исполнения. Так, для IMM выполнение инструкции  $[x]^{\text{rel}} := 1$  соответствует добавлению в текущий граф вершины с очередным порядковым номером и меткой вида  $\mathbf{W}^{\text{rel}}(x, 1)$ , а также рёбер, отражающих синтаксические зависимости данного события.

**Определение 3.4.** Граф исполнения определяет сценарий поведения программы — функцию  $f : \text{Loc} \rightarrow \text{Val}$ , отображающую адрес в последнее (согласно порядку  $\mathbf{co}$ ) записанное по нему значение.

Декларативная модель памяти задаётся предикатом согласованности, которому должны удовлетворять графы исполнения программ.

**Определение 3.5.** Сценарий поведения программы является согласованным по модели памяти  $M$ , если он задан некоторым графом её исполнения, удовлетворяющим предикату согласованности  $M$ .

## 3.4. Описание используемых моделей памяти

В данном разделе описываются рассматриваемые модели памяти —  $\text{OSamLM}$  и  $\text{IMM}$ , — а также их предикаты согласованности.

### 3.4.1. Модель памяти $\text{OSam}$ ( $\text{OSamLM}$ )

Модель памяти  $\text{OSam}$  ( $\text{OSamLM}$ ) [4] задана эквивалентными операционным и декларативным описаниями. Для доказательства корректности компиляции будет использоваться декларативное описание.

$\text{OSamLM}$  поддерживает два режима доступа: неатомарный  $\text{na}$  и атомарный  $\text{at}$  (схожи с  $\text{rl}$  и  $\text{sc}$  в  $\text{C++}$ ). При этом память также разделена на неатомарные и атомарные адреса, и к конкретному адресу можно обратиться только операцией соответствующего режима.

В графе исполнения  $\text{OSamLM}$  есть только операции чтения и записи, барьеры отсутствуют.

Перед рассмотрением предиката согласованности введём ещё несколько обозначений. Для отношения  $R$  в графе исполнения будем обозначать  $R_i$  рёбра  $R$ , проходящие между вершинами одного потока, а  $R_e$  — между вершинами разных потоков.

**Определение 3.6.** *Сценарий исполнения называется согласованным по  $\text{OSamLM}$ , если в соответствующем графе исполнения выполняются следующие аксиомы:*

- 1) *последовательная согласованность по отдельным адресам (SC per location, coherence): отношение  $\text{hbo}$ ;  $(\text{co} \cup \text{fr})$  иррефлексивно, где  $\text{hbo} \triangleq \text{po} \cup [\text{E}^{\text{at}}]; (\text{rf} \cup \text{co}); [\text{E}^{\text{at}}];$*
- 2) *отсутствие буферизации при чтении (load buffering): отношение  $\text{po} \cup \text{rfe} \cup [\text{E}^{\text{at}}]; (\text{coe} \cup \text{fre}); [\text{E}^{\text{at}}]$  ациклично.*

### 3.4.2. Промежуточная модель памяти

Промежуточная модель памяти ( $\text{IMM}$ ) определена декларативно. Полный предикат согласованности  $\text{IMM}$  достаточно сложен, поэтому мы

рассмотрим лишь часть модели, которая будет необходима для построения схемы компиляции. Перед этим введём ещё несколько обозначений.  $R_{loc}$  будем обозначать рёбра  $R$ , проходящие между вершинами с метками одного и того же адреса,  $R_{\neq loc}$  — между вершинами с метками разных адресов.

Синтаксис программ на ИММ напоминает таковой в C++ — помимо инструкций атомарного чтения и записи есть инструкции барьеров памяти, а также операций read-modify-write. В граф исполнения программы на ИММ, помимо отношений, перечисленных в разделе 3.3, также входит отношение, связывающее события чтения и записи, которые совершаются при операции read-modify-write:

$$G.rmw \subseteq ([G.R]; po|_{imm}; [G.W])_{loc}.$$

Для построения схемы компиляции мы пользуемся расширением [15] ИММ, которое дополняет оригинальную модель [18] sc-операциями.

**Определение 3.7.** *Сценарий исполнения называется согласованным по ИММ, если в соответствующем графе исполнения выполняются следующие аксиомы:*

- 1) отношение  $hb$ ;  $(rf \cup co \cup fr)^+$  иррефлексивно, при этом справедливо следующее:

$$hb \triangleq (po \cup sw)^+$$

$$sw \triangleq release; (rfi \cup po_{loc}^?; rfe); ([R^{acq}] \cup po; [F^{acq}])$$

$$release \triangleq ([W^{rel}] \cup [F^{rel}]; po); rs$$

$$rs \triangleq [W]; po_{loc}; [W] \cup [W]; (po_{loc}^?; rfe; rmw)^*;$$

- 2) операции read-modify-write являются атомарными:

$$rmw \cap (fre; coe) = \emptyset;$$

- 3) отношение  $ag$  ациклично, при этом справедливо следующее:

$$ag \supset rfe \cup bob$$

$$bob \supset [R^{acq}]; po \cup po; [F] \cup [F]; po;$$

4) отношение  $\text{psc}_{\text{base}}$  ациклично, при этом справедливо следующее:

$$\text{psc}_{\text{base}} \triangleq ([E^{\text{sc}}] \cup [F^{\text{sc}}]; \text{hb}^?); \text{scb}; ([E^{\text{sc}}] \cup \text{hb}^?; [F^{\text{sc}}])$$

$$\text{scb} \triangleq \text{po} \cup \text{po}_{\neq \text{loc}}; \text{hb}; \text{po}_{\neq \text{loc}} \cup \text{hb}_{\text{loc}} \cup \text{co} \cup \text{fr}.$$

## 4. Обоснование схемы компиляции

### 4.1. Необходимость введения дополнительных инструкций

Ранее в разделе 3.2 был рассмотрен пример программы, компиляция которой требует вставки дополнительных инструкций. Для удобства граф её исполнения приведён здесь на рис. 5. Именно, данный граф показывает, что ребро  $rf_e$  должно быть окружено барьерами  $rel$  и  $acq$ , поэтому инструкции неатомарной записи и атомарного чтения компилируются с использованием указанных барьеров.

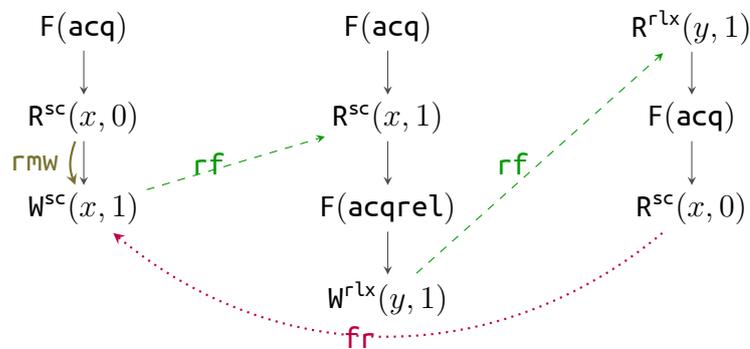


Рис. 5: Ребро  $rf_e$  между вершинами с режимом  $rlx$  должно быть окружено барьерами

Однако видно, что результирующая программа на рис. 4 получена с использованием схемы компиляции, которая вставляет в программу и другие инструкции. Так, барьер перед инструкцией неатомарной записи по адресу  $y$  имеет более строгий режим  $acqrel$ . Более того, инструкция атомарной записи компилируется в инструкцию read-modify-write, предварённую  $acq$ -барьером.

Приведённые ниже примеры демонстрируют, какие нежелательные сценарии поведения запрещают барьеры и инструкции read-modify-write. Для простоты мы рассмотрим эти случаи по отдельности, т.е. не будем вставлять барьеры в граф, содержащий read-modify-write и наоборот.

Использование барьеров позволяет обеспечить гарантируемое OSaM MM свойство multicopy atomicity для атомарных адресов. Это означает

отсутствие в графе исполнения скомпилированной программы циклов, состоящих из рёбер  $fr$ , проходящих между вершинами с меткой  $sc$ , а также рёбер  $ro$  и  $rf$ . IMM в общем случае допускает такие циклы. Поскольку отношение  $ag$  в IMM не учитывает рёбра  $fr$ , то с целью запрета подобного сценария поведения цикл вида  $[W^{at}]; (ro; rf)^+; ro; [R^{at}]; fr$  необходимо представить как  $hb^+; fr$ . Для этого каждое ребро  $rf$  между неатомарными вершинами необходимо окружить барьерами. Кроме того, так как после подобного ребра может идти  $rf$  между атомарными вершинами, перед  $W^{at}$  необходимо также расположить  $acq$ -барьер. Наконец,  $acq$ -барьер необходимо расположить и перед ребром  $fr$ , т.е. перед событием атомарного чтения. рис. 6 показывает последовательность рёбер  $rf$  между вершинами различных типов, иллюстрирующую необходимость расположения барьеров.

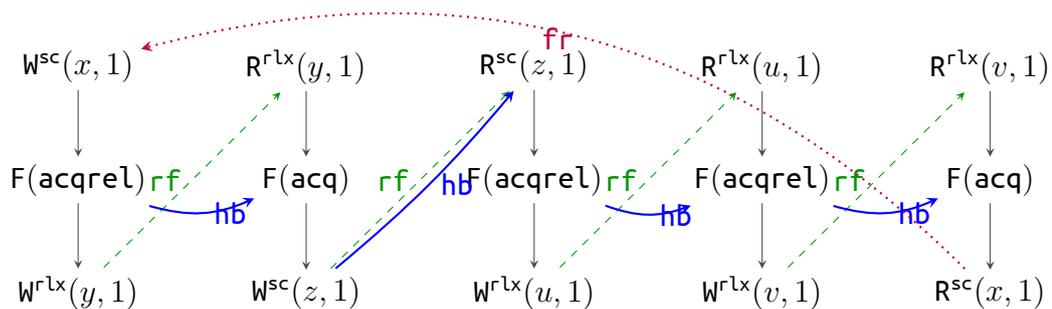


Рис. 6: Вставка барьеров позволяет провести рёбра  $hb$  вдоль рёбер  $rf$  и запретить указанный сценарий поведения

Компиляция инструкций атомарной записи в инструкции `exchange` позволяет, согласно требованиям `OSamLMM`, связать отношением  $hb$  события атомарной записи, упорядоченные по  $co$ . Вставка барьеров в этом случае не поможет, т.к. с помощью них можно провести  $hb$  лишь вдоль рёбер  $rf$ . Однако такое ребро можно ввести искусственно, добавив перед событием записи событие чтения, которое наблюдает  $co$ -предыдущую запись. Такое условие обеспечивается инструкцией `exchange`. Она порождает в графе пару вершин, связанных отношением  $rmw$ , что иллюстрируется на рис. 7.

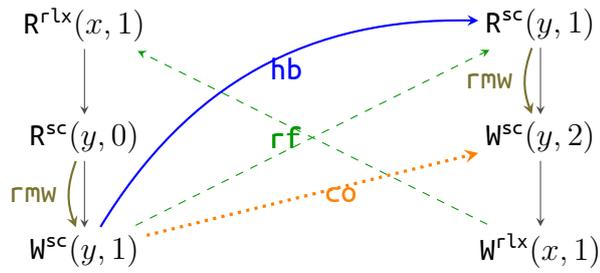


Рис. 7: Реализация атомарной записи инструкцией read-modify-write позволяет про- вести ребро **hb** и запретить указанный сценарий поведения

## 4.2. Схема компиляции

таб. 1 содержит предлагаемую схему компиляции OCamlMM в IMM. За основу взята схема компиляции OCamlMM в модель ARMv8 из [4].

Как описано в разделе 4.1, дополнительные инструкции используются для того, чтобы запретить для скомпилированной программы сценарии поведения, разрешённые IMM и запрещённые OCamlMM.

OCamlMM	IMM	ARMv8
$r := [x]^{na};$	$r := [x]^{rlx};$	$r := [x]^{rlx};$
$[x]^{na} := v;$	$\text{fence}^{\text{acqrel}}; [x]^{rlx} := v;$	$\text{fence}^{\text{acq}}; [x]^{rlx} := v;$
$r := [x]^{at};$	$\text{fence}^{\text{acq}}; r := [x]^{sc};$	$\text{fence}^{\text{acq}}; r := [x]^{sc};$
$[x]^{at} := v;$	$\text{fence}^{\text{acq}}; \text{exchange}^{sc}(x, v);$	$\text{exchange}^{sc}(x, v); \text{fence}^{\text{rel}};$

Таблица 1: Схема компиляции OCamlMM в IMM и сравнение её со схемой компиляции OCamlMM в ARMv8

## 5. Доказательство корректности компиляции

**Теорема 3.1.** Пусть  $PO = \parallel_{\tau \in \text{Tid}} PO_{\tau}$  и  $PI = \parallel_{\tau \in \text{Tid}} PI_{\tau}$  — программы для моделей  $\text{OSamMM}$  и  $\text{IMM}$  соответственно, причём для любого  $\tau \in \text{Tid}$  подпрограмма  $PI_{\tau}$  получена компиляцией  $PO_{\tau}$  с помощью схемы компиляции из таб. 1. Пусть  $G_I$  — согласованный по  $\text{IMM}$  граф исполнения  $PI$ . Тогда существует  $G_O$  — согласованный по  $\text{OSamMM}$  граф исполнения  $PO$ , соответствующий  $G_I$ .

*Доказательство.* В теореме 5.2 доказываемся, что существует  $G_O$ , являющийся графом исполнения  $PO$  и соответствующий  $G_I$  (формальное понятие соответствия графов рассматривается в разделе 5.1).

Далее, в теореме 5.1 показывается, что согласованность по  $\text{OSamMM}$   $G_O$  следует из согласованности  $G_I$  по  $\text{OSamMM}$ . Затем в теоремах 5.5 и 5.12 доказываются два условия согласованности  $G_I$  по  $\text{OSamMM}$ .  $\square$

### 5.1. Соответствие графов исполнения

Необходимо точно определить соответствие графов исполнения, чтобы один и тот же сценарий исполнения программы можно было представить как в  $\text{IMM}$ , так и в  $\text{OSamMM}$ . Отметим, что графы исполнения в этих моделях различны, т.к. в результате компиляции в программу добавляются новые инструкции, а в граф — новые вершины. Но эти различия не меняют сценарий поведения программы согласно определению 3.4.

Чтобы формализовать соответствие графов исполнения, опишем, чем отличаются множества вершин и рёбер в графах исполнения скомпилированной программы и исходной. Вершины графа исполнения в  $\text{OSamMM}$  — те же, что в графе  $\text{IMM}$ , за исключением вершин-барьеров, а также операций чтения, выполняемых в ходе операции `read-modify-write`. При удалении этих вершин также удаляются смежные им рёбра. Метки оставшихся событий в  $\text{OSamMM}$  совпадают с таковыми в  $\text{IMM}$  с точностью до переименования режимов доступа.

**Определение 5.1.** Граф исполнения по ОСамЛММ  $G_O$  соответствует графу исполнения по ИММ  $G_I$ , если выполняются следующие условия:

- 1)  $G_O.E = G_I.E \cap (G_I.R \cup G_I.W \setminus \text{dom}(G_I.rmw))$
- 2)  $\forall e. e \in G_O.E \Rightarrow G_O.\text{Lab } e = \text{renameMode}(G_I.\text{Lab } e)$ , где `renameMode` меняет метку вершины с `sc` на `at` и с `rlx` на `na`
- 3)  $G_O.\text{co} = [G_O.E]; G_I.\text{co}; [G_O.E]$
- 4)  $G_O.\text{rf} = [G_O.E]; G_I.\text{rf}; [G_O.E]$

**Теорема 5.1.** Пусть  $G_O$  и  $G_I$  — пара соответствующих графов. Тогда из согласованности  $G_I$  по ОСамЛММ следует согласованность  $G_O$  по ОСамЛММ.

*Доказательство.* Условие согласованности по ОСамЛММ требует ацикличности и иррефлексивности отношений, построенных из рёбер `ro`, `rf`, `co` и `fr`. Заметим, что  $G_O.E \subseteq G_I.E$ , так как в  $G_O.E$  отсутствуют вершины, соответствующие барьерам и операциям `read-modify-write`, а новых вершин в  $G_O$  не вводится. Кроме того, при удалении вершин из графа удаляются и смежные с ними рёбра, поэтому выполнено следующее:  $\forall r \in \{\text{ro}, \text{rf}, \text{co}, \text{fr}\}. G_O.r \subseteq G_I.r$ .

Если некоторое отношение ациклично (иррефлексивно), то и любое включённое в него отношение также ациклично (иррефлексивно). Поэтому из согласованности  $G_I$  по ОСамЛММ (при замене режимов доступа ИММ на их аналоги в ОСамЛММ) следует согласованность  $G_O$  по ОСамЛММ. □

## 5.2. Построение графа, соответствующего данному

Пусть  $PO = \parallel_{\tau \in \text{Tid}} PO_\tau$  и  $PI = \parallel_{\tau \in \text{Tid}} PI_\tau$  — программы для моделей ОСамЛММ и ИММ соответственно, причём для любого  $\tau \in \text{Tid}$  подпрограмма  $PI_\tau$  получена компиляцией  $PO_\tau$  с помощью схемы компиляции из таб. 1. Пусть  $G_I$  — согласованный по ИММ граф исполнения  $PI$ .

Следующая теорема практически повторяет теорему 3.1, за исключением требования на согласованность по ОСамЛММ искомого графа.

**Теорема 5.2.** Существует граф  $G_O$ , который задаёт сценарий исполнения  $PO$  и соответствует  $G_I$ .

*Доказательство.* Построим  $G_O$  по лемме 5.3.

Как следует из формулировки данной леммы, сужение этого графа на каждый из потоков является однопоточным графом исполнения соответствующей подпрограммы, поэтому  $G_O$  является графом исполнения  $PO$ . Осталось доказать соответствие  $G_O$  и  $G_I$ .

Условия на соответствие отношений **rf** и **co** выполняются непосредственно по формулировке леммы 5.3.

Заметим, что условия на соответствие множеств вершин  $G_O.E$  и  $G_I.E$  и функций  $G_O.Lab$  и  $G_I.Lab$  можно доказывать для сужений графов на потоки: эти условия выполняются для сужения на любой поток тогда и только тогда, когда они выполняются и во всём графе. По формулировке леммы 5.3 сужения этих графов на любой поток соответствуют друг другу.  $\square$

**Лемма 5.3.** Существует граф  $G_O$ , для которого справедливы следующие утверждения:

- $G_O.Init = G_I.Init$ ;
- $G_O.co = [G_O.E]; G_I.co; [G_O.E]$ ;
- $G_O.rf = [G_O.E]; G_I.rf; [G_O.E]$ ;
- для всех  $\tau \in Tid$  сужение  $G_{O\tau}$  соответствует  $G_{I\tau}$ , а также является графом исполнения подпрограммы  $PO_\tau$ .

*Доказательство.* Искомый граф построим как объединение однопоточных графов исполнения, полученных по лемме 5.4. К этому объединению добавим множество инициализирующих вершин  $G_I.Init$ . Отношение **co** и **rf** зададим как таковые в  $G_I$ , ограничив их на множества вершин в  $G_O$ .  $\square$

**Лемма 5.4.** Для всех  $\tau \in Tid$ ,  $PO_\tau, PI_\tau, G_{I\tau}$ , где  $G_{I\tau}$  является сужением  $G_I$  на поток  $\tau$ , существует такой граф  $G_{O\tau}$ , который является графом исполнения для программы  $PO_\tau$  и соответствует  $G_{I\tau}$ .

*Доказательство.* Сужение  $G_{I\tau}$  является однопоточным графом исполнения для соответствующего потока, так как  $G_I$  — граф исполнения для программы  $PI$ . Как упоминалось в разделе 3.3, для задания этого используется операционная семантика, которая при исполнении инструкций потока строит соответствующий граф. Состояние абстрактной машины такой семантики содержит текущий граф, значения регистров и указатель на текущую инструкцию.

Для упрощения доказательства сгруппируем переходы абстрактной машины, соответствующие исполнению блоков инструкций, каждый из которых получен компиляцией одной инструкции в  $PO_\tau$ . Мы получаем блочные состояния и переходим к т.н. блочному исполнению, в ходе которого выполнение блока переходов рассматривается как неделимый переход.

Исходное утверждение доказывается индукцией. Именно, для любого числа блочных переходов и результирующего состояния ИММ-машины можно исполнить столько же обычных переходов исполнения  $PO_\tau$  и получить состояние ОСамИММ-машины, соответствующее заданному.  $\square$

### 5.3. Доказательство последовательной согласованности по отдельным адресам

**Теорема 5.5.** *Отношение  $hb_o; (co \cup fr)$  является иррефлексивным.*

Докажем, что для выбранной схемой компиляции справедливо  $hb_o \subseteq hb$ . С учётом этого доказательство теоремы тривиально: по определению согласованности по ИММ, отношение  $hb; (rf \cup co \cup fr)$  иррефлексивно.

Для этого факта сначала покажем, что последовательность рёбер  $co$  между  $sc$ -событиями порождает  $hb$ .

**Лемма 5.6.** *Порядок событий  $sc$ -записи согласуется с отношением happens-before:  $[E^{sc}]; co; [E^{sc}] \subseteq hb$ .*

*Доказательство.* Заметим, что  $[E^{sc}]; co; [E^{sc}]$  транзитивно. Тогда можно перейти к рассмотрению непосредственных  $co$ -соседей.

Предположим, что  $[E^{sc}]; \text{co}_{imm}; [E^{sc}] \subseteq [E^{sc}]; \text{rf}; [E^{sc}]; \text{po}$ . Тогда доказательство тривиально:  $\text{rf}$  по  $sc$  событиям порождает  $\text{hb}$ , как и следующий за ним  $\text{po}$ . Значит, остаётся доказать утверждение о включении в  $[E^{sc}]; \text{rf}; [E^{sc}]; \text{po}$ .

Рассмотрим два события  $w_1, w_2 \in W^{sc}$  — непосредственных  $\text{co}$ -соседей. По схеме компиляции перед  $w_2$  следуют  $f \in F^{acq}$  и  $r \in R^{sc}$ , причём  $\text{rmw}(r, w_2)$ .

Покажем, что  $\text{rf}(w_1, r)$ . Рассмотрим событие записи  $w'$ , из которого читает  $r$ . Обращения по одному и тому же адресу имеют один и тот же режим, поэтому  $w' \in W^{sc}$ . Пусть  $w_1 \neq w'$ . Тогда либо  $\text{co}(w', w_1)$ , либо, наоборот,  $\text{co}(w_1, w')$ . В первом случае нарушается атомарность  $\text{rmw}$  между  $w_2$  и  $r$ . Во втором случае получается, что между  $\text{co}$ -соседями  $w_1$  и  $w_2$  расположен  $w'$ , что невозможно.  $\square$

**Лемма 5.7.** *Отношение happens-before в IMM содержит happens-before в  $\text{OSam} \upharpoonright \text{IMM}$ :  $\text{hbo} \subseteq \text{hb}$ .*

*Доказательство.*  $\text{hbo} \triangleq \text{po} \cup [E^{sc}]; (\text{co} \cup \text{rf}); [E^{sc}]$ . По предыдущей лемме  $\text{co}$ , ограниченный на  $sc$ , входит в  $\text{hb}$ .  $\text{rf}$  по  $sc$  событиями порождает  $\text{sw}$ , и, следовательно,  $\text{hb}$ .  $\square$

## 5.4. Доказательство отсутствия буферизации при чтении

Сначала докажем утверждения, которые позволят нам находить барьеры в программном порядке между событиями.

**Лемма 5.8.** *Перед событием записи располагается барьер:*

$$[E \setminus F]; \text{po}; [W] \subseteq \text{po}; ([F^{acqrel}]; \text{po}; [E^{rlx}] \cup [F^{acq}]; \text{po}; [E^{sc}]); [W] \cup \text{rmw}.$$

*Доказательство.* Согласно схеме компиляции, инструкция барьера располагается либо непосредственно перед инструкцией неатомарной записи, либо перед инструкцией read-modify-write. Таким образом, барьера между событием и  $\text{po}$ -следующим событием записи может не быть, только если это событие чтения в  $\text{rmw}$ .  $\square$

**Лемма 5.9.** В программном порядке между событиями  $glx$  и  $sc$  располагается барьер:  $[E^{glx}]; po; [E^{sc}] \subseteq po; [F^{acq}]; po$ .

*Доказательство.* Согласно схеме компиляции, все инструкции в режиме  $sc$  предваряются  $acq$ -барьерами.  $\square$

Кроме того, нам понадобятся следующие факты из алгебры [21].

**Лемма 5.10.** Цикл из рёбер двух типов можно представить в виде чередующихся участков рёбер каждого типа:  $(x \cup y)^+ = y^+ \cup y^*$ ;  $(x; y^*)^+$ , где  $x, y$  — произвольные отношения.

**Лемма 5.11.** Отношение  $x \cup y$  является ациклическим, если ациклически отношения  $x, y$  и  $x^+; y^+$ .

**Теорема 5.12.** Отношение  $po \cup rfe \cup [E^{sc}]; (coe \cup fre); [E^{sc}]$  является ациклическим.

Сгруппируем первые два отношения в объединении. Тогда по лемме 5.11 нужно показать ациклическость следующих отношений:

- $po \cup rfe$ ;
- $[E^{sc}]; (coe \cup fre); [E^{sc}]$ ;
- $(po \cup rfe)^+; ([E^{sc}]; (coe \cup fre); [E^{sc}])^+$ , и это эквивалентно ациклическости:  $[E^{sc}]; (po \cup rfe)^+; [E^{sc}]; ([E^{sc}]; (coe \cup fre); [E^{sc}])^+$ .

Докажем второе и третье утверждение показав, что соответствующие отношения лежат в  $([E^{sc}]; scb; [E^{sc}])^+ \subseteq psc_{base}^+$ . Напомним, что

$$scb \triangleq po \cup po_{\neq loc}; hb; po_{\neq loc} \cup hb_{loc} \cup co \cup fr \text{ и}$$

$psc_{base} \triangleq ([E^{sc}] \cup [F^{sc}]; hb^?); scb; ([E^{sc}] \cup hb^?; [F^{sc}])$ . В свою очередь, ациклическость  $psc_{base}$  следует из согласованности графа по IMM.

Теперь видно, что второе утверждение верно по определению  $scb$ . По этой же причине для доказательства третьего утверждения остаётся показать, что  $[E^{sc}]; (po \cup rfe)^+; [E^{sc}] \subseteq ([E^{sc}]; scb; [E^{sc}])^*$ .

**Теорема 5.13.** Отношение  $po \cup rfe$  ациклично.

*Доказательство.* Вновь воспользуемся леммой 5.11 и разложим условие ацикличности объединения на ацикличность отношений  $\rho$ ,  $rfe$  и  $\rho^+$ ;  $rfe^+$ . Ацикличность первого следует из согласованности по IMM, двух и более рёбер второго подряд идти не может (т.к. их концы имеют разные типы), а ацикличность третьего эквивалентна ацикличности  $\rho$ ;  $rfe$ .

Пусть такой цикл существует. Покажем, что это противоречит условию ацикличности  $ag$  (следует из согласованности по IMM). Напомним, что  $ag \supset rfe \cup bob$  и  $bob \supset [R^{acq}]; \rho \cup \rho; [F] \cup [F]; \rho$ .

По лемме 5.8 перед событием записи, которой начинается ребро  $rfe$ , есть барьер  $F^{\exists acq}$ , либо весь  $\rho$  является  $gmw$ . В первом случае внутри  $\rho$  есть барьер, а такое отношение лежит в  $bob \subseteq ag$ . Во втором случае  $gmw$  начинается с  $R^{sc}$ , и такое ребро  $\rho \supseteq gmw$  также содержится в  $bob$ . Наконец,  $rfe \subseteq ag$ .  $\square$

Теперь для доказательства второго условия согласованности по  $OCamL$  MM осталось доказать утверждение  $[E^{sc}]; (\rho \cup rfe)^+; [E^{sc}] \subseteq ([E^{sc}]; scb; [E^{sc}])^*$ .

**Теорема 5.14.** *Последовательность из рёбер  $\rho$  и  $rfe$  между вершинами  $sc$  состоит из рёбер  $scb$  между вершинами  $sc$ :  $[E^{sc}]; (\rho \cup rfe)^+; [E^{sc}] \subseteq ([E^{sc}]; scb; [E^{sc}])^*$ .*

*Доказательство.* Сначала сформулируем утверждение, которое позволит отбрасывать  $rfe$ -рёбра.

**Лемма 5.15.** *Рёбра  $rfe$ , у которых один из концов —  $sc$ , входят в  $scb$ :  $[E^{sc}]; rfe \cup [w \setminus \text{init}]; rfe; [E^{sc}] \subseteq [E^{sc}]; scb; [E^{sc}]$ .*

*Доказательство.* Если один из концов ребра  $rf$  является  $sc$ , таким же является и второй (исключение — чтение из инициализирующих записей, которые в IMM являются  $rlx$ ). Тогда  $[E^{sc}]; rf; [E^{sc}] \subseteq [E^{sc}]; hb_{loc}; [E^{sc}] \subseteq [E^{sc}]; scb; [E^{sc}]$ .  $\square$

По лемме 5.10 имеем следующее:

$[E^{sc}]; (\rho \cup rfe)^+; [E^{sc}] = [E^{sc}]; (rfe^+ \cup rfe^*; (\rho; rfe^*)^+); [E^{sc}] = [E^{sc}]; (rfe \cup rfe^?; (\rho; rfe^?)^+); [E^{sc}]$ .

Начальные участки  $rfe$ , если они есть, можно отбросить по лемме 5.15. Рассмотрим оставшееся транзитивное замыкание:

$$(po; rfe^?)^+ = po; (po; rfe)^*; po^? = po; po^? \cup po; (po; rfe)^+; po^? = po \cup (po; rfe)^+; po^?.$$

В первом случае отношение сводится к  $po \subseteq scb$ . Во втором случае, если последним ребром является  $rfe$ , то его можно отбросить по лемме 5.15. Остаётся случай  $(po; rfe)^+; po$ . К каждой паре  $po; rfe$  можно применить лемму 5.8. В результате нужно доказать следующее утверждение:

$$[(W \cup R)^{sc}]; (po; ([F^{acqrel}]; po; [E^{rlx}]; rfe \cup [F^{acq}]; po; [E^{sc}]; rfe) \cup rfw; rfe)^+; po; [(W \cup R)^{sc}] \subseteq ([E^{sc}]; scb; E^{sc})^*$$

Воспользуемся леммой 5.10. Имеем, что либо транзитивное замыкание состоит только из  $rfw; rfe$ , либо такие пары рёбер могут следовать после  $po; ([F^{acqrel}]; po; [E^{rlx}]; rfe \cup [F^{acq}]; po; [E^{sc}]; rfe)$ . В первом случае замыкание имеет вид  $hb_{loc} \subseteq scb$ , а так как оно заканчивается  $sc$ -событием, оставшееся ребро  $po$  также пройдёт по  $sc$  и образует  $scb$ . Во втором случае рассмотрим, что именно находится под транзитивным замыканием:

$$(po; ([F^{acqrel}]; po; [E^{rlx}]; rfe \cup [F^{acq}]; po; [E^{sc}]; rfe); (rfw; rfe)^*)^+ = (po; [F^{acq}]; ([F^{acqrel}]; po; [E^{rlx}] \cup [F^{acq}]; po; [E^{sc}]); (rfe; rfw)^*; rfe)^+ = (po; [F^{acq}]; C; rfe)^+, \text{ где}$$

$$C = C_1 \cup C_2 = [F^{acqrel}]; po; [E^{rlx}]; (rfe; rfw)^* \cup [F^{acq}]; po; [E^{sc}]; (rfe; rfw)^*.$$

Заметим, что верно следующее:

$$(po; [F^{acq}]; C; rfe)^+ = po; [F^{acq}]; (C; rfe; po; [F^{acq}])^*; C; rfe.$$

В результате необходимо доказать вот что:

$$[(W \cup R)^{sc}]; po; [F^{acq}]; (C; rfe; po; [F^{acq}]);^*; C; rfe; po; [(W \cup R)^{sc}] \subseteq ([E^{sc}]; scb; E^{sc})^*.$$

Заметим, что выполнено следующее утверждение:

$$(C; rfe; po; [F^{acq}])^* = (([F^{acqrel}]; po; [E^{rlx}] \cup [F^{acq}]; po; [E^{sc}]); (rfe; rfw)^*; rfe; po; [F^{acq}])^* \subseteq hb^?.$$

Это справедливо, поскольку  $hb$  задано так:

$$hb \triangleq (po \cup sw)^+, \\ sw \supset release; rfe; po; [F^{acq}],$$

$\text{release} \triangleq ([W^{\text{rel}}] \cup [F^{\text{rel}}]; \text{po}); \text{rs},$

$\text{rs} \supset (\text{rfe}; \text{rmw})^*.$

Вспомним, что  $\text{po}_{\neq \text{loc}}; \text{hb}; \text{po}_{\neq \text{loc}} \subseteq \text{scb}$ . Воспользуемся тем, что  $\text{po}$  между событием чтения/записи и барьером образует именно  $\text{po}_{\neq \text{loc}}$ . Тогда видно, что  $[(W \cup R)^{\text{sc}}]; \text{po}; [F^{\text{acq}}] \subseteq [E^{\text{sc}}]; \text{po}_{\neq \text{loc}}$ .

Остаётся доказать следующее утверждение:

$[E^{\text{sc}}]; \text{po}_{\neq \text{loc}}; \text{hb}^?; C; \text{rfe}; \text{po}; [(W \cup R)^{\text{sc}}] \subseteq ([E^{\text{sc}}]; \text{scb}; [E^{\text{sc}}])^*.$

Для этого перепишем  $C = C_1 \cup C_2$  и докажем утверждение для  $C_1$  и  $C_2$  по отдельности.

Случай с  $C_1$ . Покажем, что

$[E^{\text{sc}}]; \text{po}_{\neq \text{loc}}; \text{hb}^?; C_1; \text{rfe}; \text{po}; [(W \cup R)^{\text{sc}}] \subseteq ([E^{\text{sc}}]; \text{scb}; [E^{\text{sc}}])^*$ . Заметим, что по лемме 5.9 в последнем ребре  $\text{po}$  найдётся  $\text{acq}$ -барьер, с помощью которого можно будет построить ребро  $\text{hb}$ :

$$\begin{aligned} C_1; \text{rfe}; \text{po}; [(W \cup R)^{\text{sc}}] &= [F^{\text{acqrel}}]; \text{po}; [E^{\text{rlx}}]; (\text{rfe}; \text{rmw})^*; \text{rfe}; \text{po}; [(W \cup R)^{\text{sc}}] \\ &= [F^{\text{acqrel}}]; \text{po}; [E^{\text{rlx}}]; (\text{rfe}; \text{rmw})^*; \text{rfe}; [E^{\text{rlx}}]; \text{po}; [F^{\text{acq}}]; \text{po}; [(W \cup R)^{\text{sc}}] \\ &\subseteq \text{hb}; \text{po}_{\neq \text{loc}}; [E^{\text{sc}}]. \text{ Тогда } [E^{\text{sc}}]; \text{po}_{\neq \text{loc}}; \text{hb}^?; C_1; \text{rfe}; \text{po}; [(W \cup R)^{\text{sc}}] \\ &\subseteq [E^{\text{sc}}]; \text{po}_{\neq \text{loc}}; \text{hb}^?; \text{hb}; \text{po}_{\neq \text{loc}}; [E^{\text{sc}}] \\ &\subseteq [E^{\text{sc}}]; \text{scb}; [E^{\text{sc}}]. \end{aligned}$$

Случай с  $C_2$ . Покажем, что выполнено следующее:

$[E^{\text{sc}}]; \text{po}_{\neq \text{loc}}; \text{hb}^?; C_2; \text{rfe}; \text{po}; [(W \cup R)^{\text{sc}}] \subseteq ([E^{\text{sc}}]; \text{scb}; [E^{\text{sc}}])^*$ . Заметим, что последовательность пар рёбер  $\text{rfe}; \text{rmw}$  входит в  $\text{scb}$ :

$$\begin{aligned} C_2 &= [F^{\text{acq}}]; \text{po}; [E^{\text{sc}}]; (\text{rfe}; \text{rmw})^* \subseteq \text{po}_{\neq \text{loc}}; [E^{\text{sc}}]; ([E^{\text{sc}}]; \text{rfe}; [E^{\text{sc}}]; \text{rmw}; [E^{\text{sc}}])^*; [E^{\text{sc}}] \\ &\subseteq \text{po}_{\neq \text{loc}}; [E^{\text{sc}}]; ([E^{\text{sc}}]; \text{scb}; [E^{\text{sc}}])^*; [E^{\text{sc}}]. \text{ В этом случае имеем:} \\ [E^{\text{sc}}]; \text{po}_{\neq \text{loc}}; \text{hb}^?; C_2; \text{rfe}; \text{po}; [(W \cup R)^{\text{sc}}] \\ &\subseteq [E^{\text{sc}}]; \text{po}_{\neq \text{loc}}; \text{hb}^?; \text{po}_{\neq \text{loc}}; [E^{\text{sc}}]; ([E^{\text{sc}}]; \text{scb}; [E^{\text{sc}}])^*; [E^{\text{sc}}]; \text{rfe}; \text{po}; [(W \cup R)^{\text{sc}}] \\ &\subseteq ([E^{\text{sc}}]; \text{scb}; [E^{\text{sc}}]); ([E^{\text{sc}}]; \text{scb}; [E^{\text{sc}}])^*; ([E^{\text{sc}}]; \text{scb}; [E^{\text{sc}}])^2. \quad \square \end{aligned}$$

## 6. Формализация доказательства в Coq

В некоторых опубликованных доказательствах корректности компиляции впоследствии были найдены неточности. Например, [10] демонстрирует ошибку в схеме компиляции C++ в Power, [18] — в схеме компиляции модели Promising в Power.

Чтобы избежать подобных ошибок, доказательство из раздела 5 было решено формализовать в Coq. Это одна из систем интерактивного доказательства теорем, которые позволяют записать доказательство математического утверждения на функциональном языке программирования и автоматически проверить его корректность. В настоящее время Coq и его аналоги используются как для формализации математики [6], так и для верификации отдельных программных проектов [12].

Выбор системы Coq обусловлен наличием существующей формализации IMM и доказательств корректности компиляции для неё [17]. Данная формализация построена на основе библиотеки `hahn`, упрощающей работу с бинарными отношениями.

На данный момент формализация `OCamlMM` и доказательств корректности компиляции в IMM занимает порядка 8 тыс. строк кода на Coq. Исходный код доступен на [17].

Структура доказательства в Coq повторяет таковую в разделе 5. Формализация разбита на следующие файлы:

- `OCaml.v` — формализация `OCamlMM`;
- `ImmProgram.v` — условия на программы в `OCamlMM` и исполнения в этой модели;
- `ImmCompScheme.v` — описание схемы компиляции и её свойств;
- `ImmCompSuggest.v` — высокоуровневая структура доказательства;
- `ImmImpliesImm.v` — доказательство согласованности по `OCamlMM` графа исполнения, согласованного по IMM;

- `OmmImmSimulation.v` — доказательство отношения симуляции между исполнениями в `OCamlMM` и `IMM`;
- `GraphConstruction.v` — построение графа исполнения из однопочтовых графов исполнения;
- `BlockSteps.v` — преобразование скомпилированной программы в последовательность блоков, полученных при компиляции одной исходной инструкции;
- `CompSchemeGraph.v` — доказательство утверждения о том, что в графе исполнения скомпилированной программы выполняются предпосылки, необходимые для доказательства его согласованности по `OCamlMM`;
- `BoundedReIsProperties.v` — вспомогательные утверждения о компонентах графа, значения которых устанавливаются пошагово в ходе его построения;
- `ClosuresProperties.v` — вспомогательные утверждения про транзитивное и транзитивно-рефлексивное замыкания;
- `Utils.v` — вспомогательные утверждения про отношения;
- `ListHelpers.v` — вспомогательные утверждения про списки, в том числе те, которые доказаны в `Coq 8.10`, но не добавлены в стандартную библиотеку `Coq 8.09`, используемую в проекте.

В основе `Coq` лежит интуиционистская логика, т.е. по умолчанию в нём не используется закон исключённого третьего. Это исключает неконструктивные доказательства, что, в частности, делает возможным экстракцию доказательств на `Coq` в программы на других языках. Однако, так как экстракция данного доказательства не требуется, а применение закона исключённого третьего значительно упрощает доказательство (и в некоторых случаях является необходимым), то было решено использовать эту аксиому.

Одним из достоинств `Coq` является возможность написания пользовательских тактик — инструкций, преобразующих текущую цель доказательства. Такие тактики могут быть частично автоматизированы, т.к. включать перебор различных вариантов доказательства прозрачно от пользователя. В рамках данного проекта было разработано несколько тактик, которые значительно сократили объём доказательств. Подходящей ситуацией для их применения является перебор случаев, выкладки в которых являются схожими, но отличаются в деталях. В этом случае разумным решением является разработка специализированной тактики, которая нивелирует различия разных случаев путём автоматизации.

В рамках данного проекта применение `Coq` позволило уточнить входные условия на программы, которые до начала работы над проектом явным образом не оговаривались. Так, явным образом выражено условие на разделение атомарных и неатомарных адресов на уровне программы для ОС `am1mm`. Также от этой программы требуется выделить под инструкции `exchange` специальный регистр, который не будет использоваться в остальной части программы. Наконец, для упрощения части доказательства, касающейся блочных состояний и переходов было решено ограничить целевые адреса в инструкциях `ifgoto`, запретив переходы за пределы программы. Данное ограничение выглядит разумным, т.к. соответствующее преобразование может быть выполнено реальным компилятором.

## 7. Существующие исследования

Проблема корректности компиляции из OSaM<sub>LMM</sub> и в IMM рассматривается и в других работах. Так, авторы OSaM<sub>LMM</sub> [4] разработали для неё схемы компиляции в архитектуры x86 и ARMv8 [19]. Модель памяти x86 практически не отличается от SC [9], поэтому для корректной компиляции OSaM<sub>LMM</sub> в x86 достаточно лишь реализовать инструкцию атомарной записи с использованием `xchg`. В схеме компиляции OSaM<sub>LMM</sub> в ARMv8, в отличие от предложенной нами схемы, при компиляции инструкции неатомарной записи используется барьер  $F^{acq}$ , а не  $F^{acqrel}$ . Это объясняется тем, что в модели ARMv8 отношение `ob` (аналог `ag` в IMM) включает в себя  $rfe \cup frfe \cup coe$  по неатомарным операциям и `po` с `acq`-барьером перед событием записи, поэтому в последовательности рёбер вида  $(po; rfe)^+$  не требуется `rel`-барьер. Наконец, в [5] утверждается, что свойство локальной свободы от гонок можно реализовать в программной транзакционной памяти с использованием тех же схем компиляции в x86 и ARMv8.

В [18] приведены схемы компиляции моделей Promising [8] и RC11 [10] в IMM. Предикат консистентности RC11 требует ацикличности отношения  $po \cup rf$ , поэтому между каждым событием не-`sc` чтения и последующим событием записи необходимо располагать барьер. С этим дополнительным условием доказательство корректности компиляции значительно упрощается. В свою очередь, модель Promising является достаточно слабой, поэтому для неё корректной является тривиальная схема компиляции (при выполнении некоторых условий на инструкции `read-modify-write` в ней). Однако доказательство этого значительно сложнее, так как в нём применяется метод обхода графа исполнения, который позволяет связать операционную семантику модели Promising и декларативную семантику IMM.

Стоит отметить, что компиляция с использованием промежуточного представления программы может быть менее эффективной, чем компиляция в целевую модель памяти напрямую. Так, при композиции схем компиляции OSaM<sub>LMM</sub> в IMM 1 и IMM в ARMv8 [15] инструкция

атомарной записи OCamlMM компилируется с использованием двух барьеров, в то время как схема компиляции OCamlMM в ARMv8 из [4] использует только один.

## 8. Заключение

В данной работе представлена схема компиляции OCamLMM в промежуточную модель, позволяющая получить схему компиляции OCamLMM в Power. При этом были получены следующие результаты.

- Предложена схема компиляции OCamLMM в IMM, использующая барьеры памяти и инструкции read-modify-write для необходимой синхронизации доступов к памяти.
- Была доказана корректность предложенной схемы. Для этого было введено формальное понятие соответствия графов, которое позволило связать сценарии исполнения исходной и скомпилированной программ. Затем было доказано, что для каждого согласованного по IMM графа исполнения скомпилированной программы существует соответствующий ему граф исполнения исходной программы, согласованный по OCamLMM.
- Полученное доказательство было формализовано в Coq. Без учёта существующей формализации IMM формализация заняла порядка 8 тыс. строк кода.

Результаты работы были представлены на Открытой конференции ИСП РАН и опубликованы в "Трудах Института системного программирования РАН" [23].

Данное исследование может быть продолжено при реализации свойства локальной свободы от гонок в других моделях памяти. Например, изучается возможность реализовать это свойство в модели памяти C++ [3]. Результаты, полученные в данной работе, могут быть использованы для построения схем компиляции таких моделей. В частности, сравнение схемы компиляции OCamLMM в ARMv8 [4] и полученной в данной работе позволяет предположить, что с большой вероятностью для компиляции инструкции атомарной записи в подобных моделях необходимо будет использовать инструкцию атомарной замены.

## Список литературы

- [1] Jade Alglave, Luc Maranget, and Michael Tautschnig. “Herding Cats: Modelling, Simulation, Testing, and Data Mining for Weak Memory”. In: *ACM Trans. Program. Lang. Syst.* 36.2 (July 2014), 7:1–7:74. ISSN: 0164-0925. DOI: [10 . 1145 / 2627752](https://doi.org/10.1145/2627752). URL: <http://doi.acm.org/10.1145/2627752>.
- [2] Mark Batty et al. “Mathematizing C++ Concurrency”. In: *POPL 2011*. ACM, 2011, pp. 55–66. DOI: [10 . 1145 / 1925844 . 1926394](https://doi.org/10.1145/1925844.1926394).
- [3] Simon Doherty. *Local Data-race Freedom and the C11 Memory Model*. Surrey Concurrency Workshop. 2019.
- [4] Stephen Dolan, KC Sivaramakrishnan, and Anil Madhavapeddy. “Bounding Data Races in Space and Time”. In: *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation*. 2018.
- [5] Brijesh Dongol, Radha Jagadeesan, and James Riely. “Modular Transactions: Bounding Mixed Races in Space and Time”. In: *Proceedings of the 24th Symposium on Principles and Practice of Parallel Programming*. PPOPP ’19. Washington, District of Columbia: Association for Computing Machinery, 2019, pp. 82–93. ISBN: 9781450362252. DOI: [10 . 1145 / 3293883 . 3295708](https://doi.org/10.1145/3293883.3295708). URL: <https://doi.org/10.1145/3293883.3295708>.
- [6] Georges Gonthier. “The Four Colour Theorem: Engineering of a Formal Proof”. In: *Computer Mathematics: 8th Asian Symposium, ASCM 2007, Singapore, December 15-17, 2007. Revised and Invited Papers*. Berlin, Heidelberg: Springer-Verlag, 2008, p. 333. ISBN: 9783540878261. URL: [https://doi.org/10.1007/978-3-540-87827-8\\_28](https://doi.org/10.1007/978-3-540-87827-8_28).
- [7] IBM Power Systems. *IBM Power Systems Facts and Features: Enterprise and Scale-out Systems with POWER8 Processor Technology*. <https://www.ibm.com/downloads/cas/JDRZDG0A>. 2018.
- [8] Jeehoon Kang et al. “A Promising Semantics for Relaxed-memory Concurrency”. In: *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages*. POPL 2017. Paris, France: ACM, 2017, pp. 175–189. ISBN: 978-1-4503-4660-3. DOI: [10 . 1145 / 3009837 . 3009850](https://doi.org/10.1145/3009837.3009850). URL: <http://doi.acm.org/10.1145/3009837.3009850>.
- [9] Ori Lahav and Viktor Vafeiadis. “Explaining Relaxed Memory Models with Program Transformations”. In: *FM 2016: Formal Methods*. Ed. by John Fitzgerald et al. Cham: Springer International Publishing, 2016, pp. 479–495. ISBN: 978-3-319-48989-6.

- [10] Ori Lahav et al. “Repairing Sequential Consistency in C/C++11”. In: *SIGPLAN Not.* 52.6 (June 2017), pp. 618–632. ISSN: 0362-1340. DOI: [10.1145/3140587.3062352](https://doi.org/10.1145/3140587.3062352). URL: <http://doi.acm.org/10.1145/3140587.3062352>.
- [11] Leslie Lamport. “How to Make a Multiprocessor Computer That Correctly Executes Multiprocess Programs”. In: *IEEE Transactions on Computers* C-28 9 (1979), pp. 690–691. URL: <https://www.microsoft.com/en-us/research/publication/make-multiprocessor-computer-correctly-executes-multiprocess-programs/>.
- [12] Xavier Leroy. “Formal verification of a realistic compiler”. In: *Communications of the ACM* 52.7 (2009), pp. 107–115. URL: <http://xavierleroy.org/publi/compcert-CACM.pdf>.
- [13] Jeremy Manson, William Pugh, and Sarita V. Adve. “The Java Memory Model”. In: *POPL 2005*. ACM, 2005, pp. 378–391. DOI: [10.1145/1040305.1040336](https://doi.org/10.1145/1040305.1040336).
- [14] Batty Mark et al. *C/C++11 mappings to processors*. <https://www.cl.cam.ac.uk/~pes20/cpp/cpp0xmappings.html>. 2016.
- [15] Evgenii Moiseenko et al. “Reconciling Event Structures with Modern Multiprocessors”. In: (Nov. 2019), Accepted for European Conference on Object–Oriented Programming 2020.
- [16] Scott Owens, Susmit Sarkar, and Peter Sewell. “A Better x86 Memory Model: x86-TSO”. In: *Theorem Proving in Higher Order Logics*. Ed. by Stefan Berghofer et al. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 391–407. ISBN: 978-3-642-03359-9.
- [17] Anton Podkopaev, Ori Lahav, and Egor Namakonov. *Intermediate Memory Model*. <https://github.com/weakmemory/imm>. 2019.
- [18] Anton Podkopaev, Ori Lahav, and Viktor Vafeiadis. “Bridging the Gap Between Programming Languages and Hardware Weak Memory Models”. In: *Proc. ACM Program. Lang.* 3.POPL (Jan. 2019), 69:1–69:31. ISSN: 2475-1421. DOI: [10.1145/3290382](https://doi.org/10.1145/3290382). URL: <http://doi.acm.org/10.1145/3290382>.
- [19] Christopher Pulte et al. “Simplifying ARM Concurrency: Multicopy-atomic Axiomatic and Operational Models for ARMv8”. In: *Proc. ACM Program. Lang.* 2.POPL (Dec. 2017), 19:1–19:29. ISSN: 2475-1421. DOI: [10.1145/3158107](https://doi.org/10.1145/3158107). URL: <http://doi.acm.org/10.1145/3158107>.
- [20] Coq development team. *The Coq Proof Assistant*. <https://coq.inria.fr/>.
- [21] Viktor Vafeiadis, Ori Lahav, and Anton Podkopaev. *Hahn library*. <https://github.com/vafeiadis/hahn/>. 2020.

- [22] Conrad Watt et al. “Repairing and mechanising the JavaScript relaxed memory model”. In: *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation* (2020). DOI: [10 . 1145 / 3385412 . 3385973](https://doi.org/10.1145/3385412.3385973). URL: <http://dx.doi.org/10.1145/3385412.3385973>.
- [23] Егор Намаконов и Антон Подкопаев. «Компиляция модели памяти OCaml в Power». В: *Труды ИСП РАН* 31.5 (2019), с. 63—78. DOI: [10 . 15514 / ISPRAS - 2019 - 31\(5\) - 4](https://doi.org/10.15514/ISPRAS-2019-31(5)-4).