

Санкт-Петербургский государственный университет

Кузнецов Илья Александрович

Выпускная квалификационная работа

Возможности среды разработки для
анализа, модификации и ассистирования
при написании кода на Python

Уровень образования: магистратура

Направление *09.03.04 «Программная инженерия»*

Основная образовательная программа *СВ.????.???? «Программная инженерия»*

Научный руководитель:
профессор кафедры Системного программирования, д.т.н. Д. В. Кознов

Консультант:
ведущий инженер ключевых проектов ООО «Техкомпания Хуавэй» Н. В. Тропин

Рецензент:
ведущий инженер ключевых проектов ООО «Техкомпания Хуавэй» ?. ?. ?

Санкт-Петербург
2024

Saint Petersburg State University

Ilya Alexandrovich Kuznetsov

Master's Thesis

Code Analysis, Modification and Assistance Python IDE Features

Education level: master

Speciality *09.03.04 «Software Engineering»*

Programme *CB.???????? «Software Engineering»*

Scientific supervisor:
professor System Programming's chair, D.E.Sc. D. V. Koznov

Consultants:
principal engineer at «Huawei Technologies CO.LTD» N. V. Tropin

Reviewer:
principal engineer at «Huawei Technologies CO.LTD» ?. ?. ?

Saint Petersburg
2024

Оглавление

Введение	4
1. Постановка цели и задач	6
2. Обзор предметной области	7
2.1. Инструменты разработки для Python	7
2.2. Возможности инструментов разработки для Python	7
2.3. SRC IDE	8
2.4. IntelliJ Platform	13
2.5. Сравнение моделей кода	16
Заключение	18
Список литературы	19

Введение

История программирования неразрывно связана с развитием инструментов разработки. От функциональности и удобства использования интегрированных сред разработки (IDE) зависит скорость создания программного обеспечения и его качество. В связи с определяющим влиянием таких средств разработки на процессы технологических компаний, сфера разработки IDE является перспективным направлением бизнеса.

Язык Python применяется для решения широкого спектра задач, возникающих в различных сферах деятельности бизнеса, в том числе для аналитики, обработки больших данных и машинного обучения. Являясь самым популярным языком программирования в 2024 году по версии IEEE [2], он создает необходимость поддерживать его во всех конкурентоспособных IDE. Язык Python является динамически строго типизированным, что, с одной стороны, позволяет быстро разрабатывать на нем прикладные программы, но, с другой стороны, оставляет простор для допущения неявных ошибок. Исходя из этого, создание продвинутых средств разработки для языка Python является актуальной, однако непростой задачей.

В рамках Saint-Petersburg Research Center (SRC), входящего в состав известной технологической компании¹, было решено создать собственную среду разработки для языка Python, основой которой является мультязыковая SRC IDE. На момент начала работы данная среда разработки уже была доступна для пользователей и поддерживала язык Java. Для языка Python велась подготовка к выпуску нового продукта в семействе SRC IDE [1]: в разработке находилась модель кода, готовность которой позволяла приступить к реализации основных возможностей языкового сервера и клиента, но в то же время она требовала апробации, доработки и стабилизации базовых интерфейсов.

Таким образом, необходимо с использованием модели кода реализовать основные возможности SRC IDE для анализа, модификации и

¹Название компании не указано с целью соблюдения соглашения о неразглашении (NDA).

помощи при написании кода на Python. К разработанному решению требуется создать инфраструктуру для тестирования и мониторинга, а также выполнить апробацию с помощью мониторинга и пользователей SRC IDE для Python.

1. Постановка цели и задач

Целью выпускной квалификационной работы является реализация с использованием модели кода основных возможностей SRC IDE для Python: анализ кода, действий для исправления проблем и реструктуризации кода, а также ассистентов для написания кода.

Для достижения цели выпускной квалификационной работы были поставлены следующие задачи.

1. Провести обзор моделей кода и возможностей в инструментах разработки для языка Python.
2. Реализовать основные возможности SRC IDE для языка Python с использованием модели кода.
3. Провести многостороннее тестирование разработанного решения, создав соответствующую инфраструктуру.
4. Настроить многосторонний мониторинг разработанного решения, создав соответствующую инфраструктуру.
5. Выполнить апробацию разработанного решения с помощью мониторинга и пользователей SRC IDE для Python.

2. Обзор предметной области

В данной главе будут рассмотрены инструменты разработки для Python, их возможности и модели кода, с помощью которых реализуется данная функциональность.

2.1. Инструменты разработки для Python

Для языка Python существует множество различных инструментов разработки, наиболее популярными являются статические анализаторы и среды разработки.

- Статические анализаторы – MyPy и PyRight – предлагают долгий, но наиболее точный анализ кода на Python. Их можно интегрировать в IDE, и они часто используются для оценки качества кода.
- Среды разработки – SRC IDE, PyCharm и VSCode (Pylance) – кроме анализа кода, который должен быть быстрым и как можно более точным, предлагают исправления проблем и действия по реструктуризации кода, а также ассистируют при написании кода на Python. И это лишь базовая часть их возможностей, поэтому они часто используются в продуктовой разработке.

2.2. Возможности инструментов разработки для Python

Теперь рассмотрим некоторые возможности, которые предоставляют инструменты разработки.

- Анализы (в PyCharm – инспекции, в VSCode – диагностики) – обнаруживают и интерпретируют проблемы в коде, определяя исправления.
- Исправления (фиксы) – это действия по модификации кода для устранения проблем, найденных инспекциями.

- Реструктуризации (рефакторинг) – это действия по модификации кода для изменения его структуры, абстракции или упрощения.
- Ассистенты (помощники) – представляют информацию о коде в удобном виде, упрощая его понимание и снижая вероятность ошибок.

Производительность и точность возможностей инструментов разработки напрямую зависят от модели кода. А их качество и полнота определяется соответствием спецификациям Python и библиотек, покрытием различных сценариев как в правильном, так и в неправильном коде, а также проработкой пользовательского опыта и дизайном.

2.3. SRC IDE

SRC IDE — это мультязыковая среда разработки, первоначально создававшаяся для внутренних нужд известной технологической компании, но подлежащая выпуску в виде продукта в будущем. На момент начала работы она поддерживала работу только с одним языком — Java. Но теперь на основе существующего решения предстоит создание семейства продуктов, чтобы обеспечить поддержку Python в обновлённой среде разработки. Её разработку планируется выполнять в общем технологическом процессе с решением для языка Java и путём создания повторно используемой инфраструктуры в текущей архитектуре, с сохранением ключевых технологий проекта. Перейдём к их рассмотрению.

2.3.1. Архитектура

Архитектурно SRC IDE состоит из двух основных компонентов: языкового сервера (Language Server) и клиентской части на основе Visual Studio Code [?] с собственным пользовательским интерфейсом, соединённых по протоколу языкового сервера (Language Server Protocol) [?], основанном на технологии JSON-RPC. Возможности IDE могут быть

гибко изменены при помощи расширений (Extension) для клиентской части IDE. В составе таких расширений могут подключаться дополнительные языковые сервера.

Языковой сервер основывается на кодовой модели, которой делегирует обработку пользовательского кода в виде проектов (Project). Многие компоненты языкового сервера основываются на интерфейсах кодовой модели для обеспечения как базовых, так и продвинутых возможностей IDE. К таким можно отнести: структуру проекта (Project Structure), структуру файла (File Structure), навигацию по коду (Navigation), текстовый поиск (Text Search), поиск использований кода (Find Usages), всплывающую информацию о коде (Hover), информацию о сигнатурах (Signature Help), рефакторинг кода (Refactoring), быстрые исправления (Quick Fixes), проверки кода (Inspection), автодополнение (Completion) и другие. Пример работы языкового сервера по протоколу LSP изображён на рисунке 1.

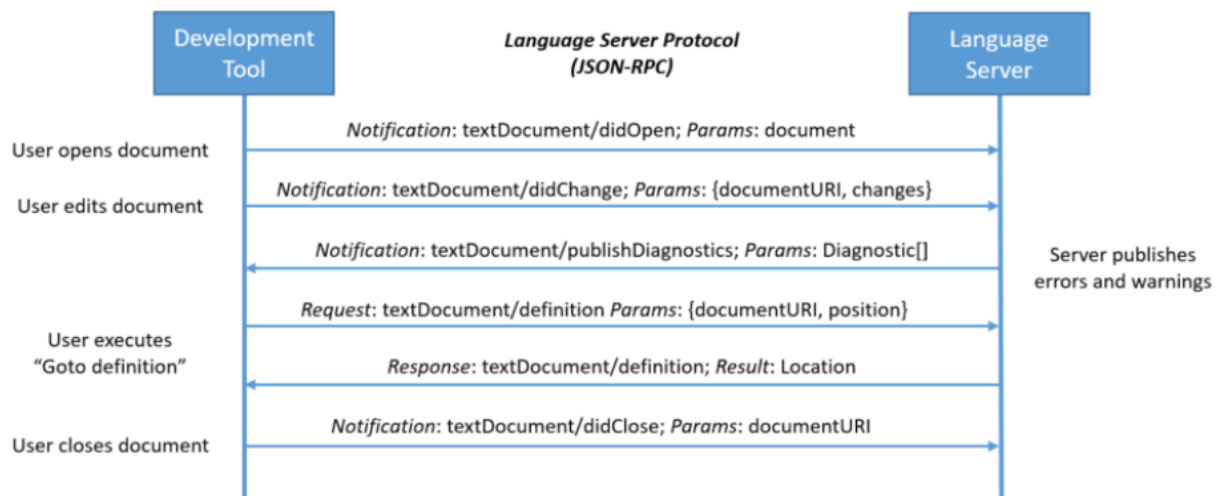


Рис. 1: Диаграмма последовательности сценария работы LS по протоколу LSP [?]

Кодовая модель отражает детальное представление пользовательского проекта в IDE с учетом структуризации кода, задаваемой языком программирования, пакетным менеджером или системой сборки (пакеты, модули). Она осуществляет разбор исходного кода и генерацию синтаксических и семантических деревьев разной степени детализации,

а также предоставляет возможности их модификации. На рисунке 2 на уровне компонентов разбирается устройство SRC IDE для языка Java.

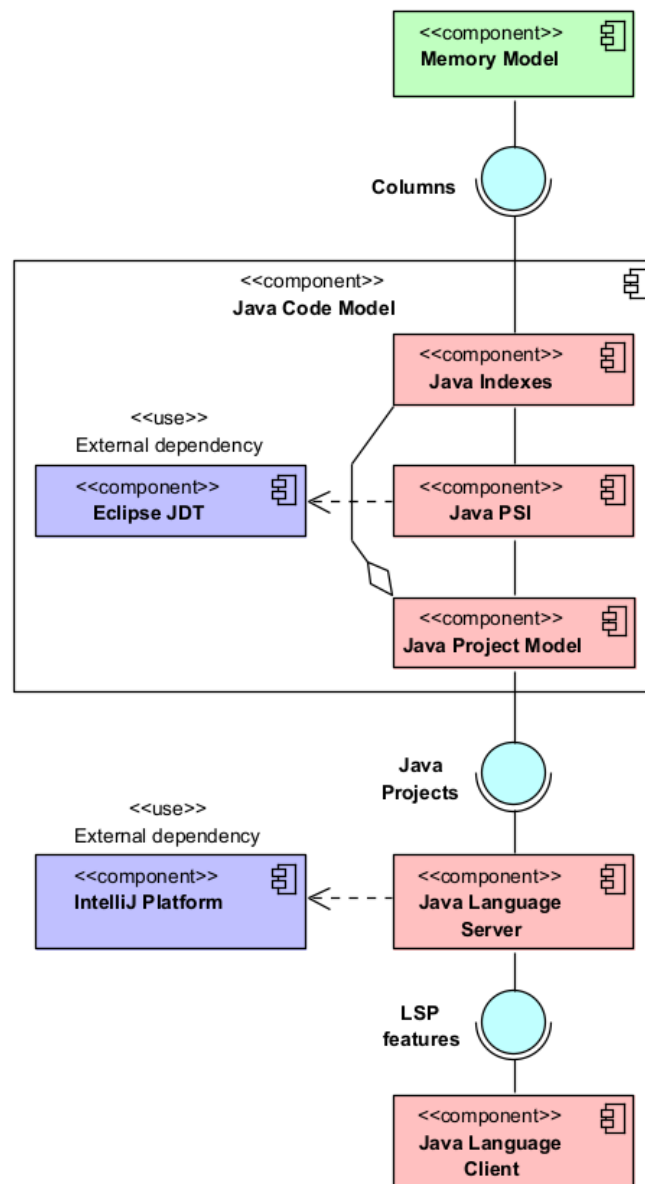


Рис. 2: Диаграмма компонентов UML для существующего решения для Java в SRC IDE

За построение индексов (Indexes) в SRC IDE отвечает собственная модель памяти (Memory Model), представленная в виде системы сжатия объектов в памяти (Object Compression System). Она предоставляет интерфейсы отображений (Mapping), которые используются для компактного хранения и быстрого доступа к классам, необходимым для работы большинства возможностей среды разработки.

Рассмотрим подробнее основные компоненты кодовой модели SRC IDE.

2.3.2. Проектная модель

Модель кода включает в себя проектную модель (Project Model). Это высокоуровневый компонент системы, отвечающий за структуру проекта и соответствующих ему зависимостей с точки зрения файловой системы, обеспечивающий конфигурирование и интеграцию с различными пакетными менеджерами и системами сборки.

В нём строится представление всего проекта в виде древовидной иерархии. Так, в отдельных логических контейнерах хранятся файлы с пользовательским кодом и зависимостями, а также их содержимое, которые распределены в структуре по соответствующим модулям и пакетам. На уровне проекта происходит объединение с возможностями других подсистем кодовой модели. Таким образом, интерфейсы проектной модели являются основными для потребления остальными частями SRC IDE. При этом хранение значительной части данных происходит в индексах.

2.3.3. Синтаксическая структура

Данный компонент позволяет по грамматике языка строить абстрактные синтаксические деревья (AST) разной степени детализации, то есть преобразовывать конкретный синтаксис в удобную для дальнейшей обработки форму. Таким образом, здесь происходит разбор файлов с кодом, содержащихся в проекте, и предоставляется доступ к соответствующим синтаксическим деревьям в виде индексов.

2.3.4. Семантическая структура

Данная подсистема осуществляет семантический анализ построенных синтаксических деревьев обнаруженных файлов с кодом. Здесь центральной сущностью является символ (Symbol) — высокоуровневый

элемент (Member) языка, являющийся абстракцией над узлами дерева. С его помощью можно устанавливать семантические связи между различными файлами с кодом, обрабатывать их зависимости, а также осуществлять разрешение простых имён (Simple Resolve) и квалифицированных имён (Qualified Resolve). Для этого также вводится такое понятие как пространство имён (Namespace), отвечающее за доступность символов в конкретном контексте, на основе которого задаётся область видимости. Доступ к информации также осуществляется при помощи индексов.

2.3.5. Модифицируемая структура

Данная подсистема отвечает за модификацию деревьев в ответ на изменения в коде и сильно связана с предыдущими двумя компонентами. Более подробно она описана в работе [?].

2.3.6. Модель памяти

В SRC IDE для управления памятью и её оптимизации, хранения данных в индексах и на диске используется система сжатия объектов. Её особенностью является компактное (сжатое) представление объектов, которое не только снижает необходимое для функционирования среды разработки количество оперативной памяти, но и обеспечивает прирост производительности благодаря увеличению количества объектов, помещающихся в кэше процессора. Таким образом, максимизируется число кэш-попаданий, за счёт чего обращаться к оперативной памяти приходится реже.

Модель памяти позволяет упаковать объекты и их поля в колонки (Column). С точки зрения времени компиляции, разницы между сохраняемыми и возвращаемыми объектами из системы сжатия не существует. Однако во время исполнения модель памяти на основе сохранённых объектов генерирует прокси-классы с точно такими же программными интерфейсами. Таким образом, при доступе к системе возвращается компактный объект прокси-класса, соответствующий исходному клас-

су. Также для работы с некоторыми сжимаемыми объектами и управления их жизненным циклом применяется технология JavaBeans.

Помимо указанных преимуществ, система сжатия имеет ограничения: граф сохраняемых объектов должен заканчиваться исключительно примитивными типами, не допускаются циклы (из-за чего иногда приходится создавать обёртки над классами), также устанавливаются определённые правила именования элементов сохраняемых классов (для внутренней работы с рефлексией). Для конфигурирования системы сжатия и упаковываемых объектов предоставляется API с аннотациями.

Для построения индексов модель памяти предоставляет гибкий интерфейс с отображениями (Mapping) разных видов: один-к-одному (OneToOne) и один-ко-многим (OneToMany). Они определяют количество значений, которые могут соответствовать одному ключу. В дополнение, предоставляются оптимизированные для работы IDE реализации отображений.

2.3.7. Индексы

Индексы являются основным компонентом системы для сохранения данных в модели памяти SRC IDE и последующего доступа к ним. Они основываются на колонках модели памяти. Остальные подсистемы как кодовой модели, так и других частей среды разработки используют существующие или собственные индексы.

Индексы образуют ациклический граф, задающий зависимость по данным одних индексов от других. Это важно для их корректного построения и обновления в ответ на поступающие изменения в виде разницы в данных (Delta). Также благодаря им осуществляется кэширование часто используемых данных, сохранение на диск и загрузка с него.

2.4. IntelliJ Platform

IntelliJ IDEA от компании JetBrains является одной из ведущих IDE, получившей немалую популярность среди разработчиков. Важным пре-

имуществом является наличие поддержки в одном продукте не только JVM-ориентированных языков Java, Kotlin и Scala, но и множества других, поставляемых в виде плагинов (Plugin). К таким можно отнести Python, Ruby, TypeScript и прочие. Альтернативным предложением компании являются производные продукты из их семейства IDE, специализированные для разработки с определённым набором технологий: PyCharm, RubyMine, WebStorm и другие [?]. Такое разнообразие решений возможно поддерживать благодаря грамотно созданной повторно используемой инфраструктуре, образующей мультязыковую IntelliJ Platform [?].

2.4.1. Архитектура

IntelliJ Platform обладает модульной архитектурой, скрывающей от разработчиков детали реализации кодовой модели и предоставляющей им набор инструментов (SDK) для расширения возможностей платформы путём добавления плагинов. Возможные точки расширения (Extension Points) определяются спецификацией платформы. Различные продукты компании собираются из платформы и включенных в поставку расширений.

Для поддержания модульной структуры, платформа разделяется на иерархию компонентов разного уровня: всего приложения (Application), проекта (Project) и модуля (Module), жизненным циклом которых управляют контейнеры, поддерживающие инъекцию зависимостей (Dependency Injection). Помимо этого, для оптимизации производительности существуют сервисы (Service), поддерживающие ленивую инициализацию.

2.4.2. Виртуальная файловая система

Работа с нативной файловой системой организована через абстракцию — виртуальную файловую систему (Virtual File System). Это позволяет платформе работать с различными файловыми системами по одному интерфейсу, обрабатывать изменения от них, а свойство перси-

стентности, за счёт сохранения снимка (Snapshot) всех файлов проекта, решает проблемы с синхронизацией доступа.

2.4.3. Инфраструктура обмена сообщениями

Для взаимодействия различных компонентов платформы создана инфраструктура обмена сообщениями, основанная на паттерне «издатель-подписчик». Это позволяет определять иерархию распространения сообщений в системе и ограничивать круг их получателей. Таким образом, существуют специализированные каналы сообщений, называемые «топики» (Topic), к которым может подключиться подписчик и ожидать получение сигналов. Издатель может посылать сообщения в соответствующую шину (Bus), которые будут впоследствии распространены всем ожидающим подписчикам.

2.4.4. Проектная модель

Проектная модель IntelliJ Platform мультязычна и расширяема. Она определяет общий интерфейс проекта, инкапсулирующего весь исходный код и поддерживающего разделение на модули, а также зависимые библиотеки и SDK. Для определения пользовательского кода применяются логические контент-руты (Content Root), позволяющие распознать, где в проекте лежит исходный код, код зависимостей и тестов, а какой следует исключить из рассмотрения. Они же влияют и на способ обработки кода в среде разработки. Также в проектной модели обеспечивается интеграция с различными системами сборки.

В связи с монолитной архитектурой IntelliJ Platform многие компоненты проектной и кодовой моделей содержат методы для реагирования на события как других подсистем, так и пользовательского интерфейса. Подобное взаимодействие происходит через инфраструктуру обмена сообщениями. Например, IntelliJ Platform предоставляет возможности для работы с пользовательским интерфейсом и используется в качестве клиентской части в среде разработки для .NET — Rider. При этом вся функциональность кодовой модели вынесена в отдельный

процесс языкового сервера на .NET, использующий собственный протокол RD (Reactive Distributed) для взаимодействия с платформой.

2.4.5. Кодовая модель

Роль модели кода в IntelliJ Platform играет мультязыковой интерфейс программной структуры (PSI). Он отвечает за синтаксический анализ файлов и создание синтаксической и семантической модели кода в виде деревьев, обеспечивающей работу многих функций платформы. После обработки пользовательского кода, его элементы доступны через PSI-файлы как PSI-элементы, с которыми удобно работать для реализации продвинутых возможностей платформы. Необходимо отметить, что интерфейс программной структуры — модифицируемый. Изменения приходят от виртуальной файловой системы и обновляют синтаксические деревья, за чем следует обновление изменившейся информации в PSI. При этом для хранения данных и поддержания модели кода в производительном состоянии используются индексы.

2.4.6. Индексы

Быстрый доступ к хранимым данным платформы реализуется на основе индексов, которые строятся и поддерживаются в актуальном состоянии на протяжении всей работы платформы. Они основываются на произвольных отображениях ключей в значения и хранятся в компактном формате. Индексы активно применяются во всех компонентах системы, в том числе для обеспечения общего доступа к данным между ними.

2.5. Сравнение моделей кода

Исходные характеристики модели кода могут быть оценены различными метриками. В статьях часто встречается метрика — количество разрешённых квалифицированных ссылок, поскольку для этого задействуются многие подсистемы модели кода, например Resolve и Type Inference.

В таблице представлено сравнение с другими решениями на основе этой метрики на конец 2023 года, новые данные ещё собираются. SRC IDE обходила все существующие решения на популярных проектах, а PyLance даже не останавливался на некоторых из них. Таким образом, при даже при идентичной реализации, например, какого-нибудь анализа, ожидается, что в SRC IDE он будет точнее.

	PyLance (PyRight)	PyCharm	SRC IDE
Flask (5072 QR)	32.96 %	29.20 %	39.33 %
Jedi (10884 QR)	32.77 %	45.61 %	60.63 %
PyTorch (599914 QR)		54.29 %	61.36 %
Tensorflow (491977 QR)		65.96 %	84.09 %
Django (202578 QR)	28.75 %	59.16 %	76.08 %
Numpy (85004 QR)		61.51 %	87.64 %
ToolBench (2913 QR)	57.56 %	62.92 %	79.71 %
SciPy (123591 QR)		32.37 %	85.63 %
Scikit-learn (90175 QR)		33.78 %	69.52 %

Рис. 3: Сравнение разрешенных квалифицированных ссылок в моделях кода

Заключение

Для достижения цели выпускной квалификационной работы были получены следующие результаты.

- Рассмотрены модели кода и возможности в статических анализаторах — MyPy, PyRight, и средах разработки — SRC IDE, PyCharm, VSCode (Pylance).
- На основе модели кода были реализованы основные возможности SRC IDE для Python, а именно: анализы кода, действия для исправления проблем и реструктуризации кода, а также ассистенты для написания кода.
- Проведено модульное, интеграционное, регрессионное и комплексное тестирование разработанного решения, создана соответствующая инфраструктура.
- Настроен мониторинг производительности и точности разработанного решения, создана инфраструктура на основе базы данных VictoriaMetrics и системы визуализации Grafana.
- Выполнена апробация разработанного решения на основе мониторинга, а также отзывов, сценариев и обнаруженных проблем у пользователей SRC IDE для Python.

Список литературы

- [1] Кузнецов И. А. Повторно используемая инфраструктура мультязыковой среды разработки для языка Python. — 2023.
- [2] Рейтинг языков программирования. — Редакция Spectrum IEEE, 2024. — URL: <https://spectrum.ieee.org/top-programming-languages-2024> (дата обращения: 15 декабря 2024 г.).