

# Оглавление

---

1 ■ Введение в системы глубокого обучения .....	27
2 ■ Служба управления наборами данных .....	63
3 ■ Служба тренировки моделей.....	120
4 ■ Распределенная тренировка .....	160
5 ■ Служба гиперпараметрической оптимизации .....	199
6 ■ Конструирование службы раздачи моделей.....	233
7 ■ Раздача моделей на практике .....	260
8 ■ Склад метаданных и артифактов.....	320
9 ■ Оркестровка рабочего процесса .....	343
10 ■ Путь к производству.....	375

# Содержание

---

<i>Оглавление</i> .....	5
<i>Вводное слово</i> .....	12
<i>Предисловие</i> .....	14
<i>Благодарности</i> .....	17
<i>О книге</i> .....	20
<i>Об авторах</i> .....	25
<i>Об иллюстрации на обложке</i> .....	26
<b>1      Введение в системы глубокого обучения</b> .....	27
1.1    Цикл освоения глубокого обучения .....	30
1.1.1    Фазы цикла освоения продукта на базе глубокого обучения ....	32
1.1.2    Технические роли в цикле освоения .....	38
1.1.3    Пошаговый обход цикла освоения глубокого обучения .....	41
1.1.4    Масштабируемость при разработке проекта.....	43
1.2    Обзор конструкции системы глубокого обучения.....	44
1.2.1    Эталонная системная архитектура.....	45
1.2.2    Ключевые компоненты .....	46
1.2.3    Ключевые пользовательские сценарии .....	52
1.2.4    Выведение своей собственной конструкции .....	55
1.2.5    Разработка компонентов поверх Kubernetes .....	57
1.3    Разработка системы глубокого обучения в сравнении с разработкой модели .....	60
Резюме .....	61
<b>2      Служба управления наборами данных</b> .....	63
2.1    Понимание службы управления наборами данных .....	65
2.1.1    Почему системы глубокого обучения нуждаются в управлении наборами данных.....	65
2.1.2    Принципы конструирования службы управления наборами данных.....	70

2.1.3	Парадоксальный характер наборов данных .....	73
2.2	Экскурсия по образцу службы управления наборами данных .....	75
2.2.1	Ознакомление с образцом службы .....	75
2.2.2	Пользователи, пользовательские сценарии и общая картина.....	82
2.2.3	API приема данных .....	84
2.2.4	API доставки тренировочного набора данных.....	90
2.2.5	Внутреннее хранилище наборов данных.....	96
2.2.6	Схемы данных.....	99
2.2.7	Добавление нового типа набора данных ( <i>IMAGE_CLASS</i> ) ....	103
2.2.8	Резюме конструкции службы .....	104
2.3	Подходы с открытым исходным кодом.....	105
2.3.1	<i>Delta Lake</i> и <i>Petastorm</i> с семейством Apache Spark.....	105
2.3.2	<i>Pachyderm</i> с облачным хранилищем объектов.....	113
	Резюме .....	118

<b>3</b>	<b>Служба тренировки моделей.....</b>	120
3.1	Служба тренировки моделей: обзор конструкции .....	122
3.1.1	Зачем использовать службу тренировки моделей? .....	123
3.1.2	Принципы конструирования службы тренировки .....	125
3.2	Шаблон исходного кода тренировки для глубокого обучения .....	127
3.2.1	Рабочий процесс тренировки моделей .....	127
3.2.2	Докеризация исходного кода тренировки моделей в качестве черного ящика.....	129
3.3	Образец службы тренировки моделей .....	130
3.3.1	Ознакомление со службой .....	131
3.3.2	Обзор конструкции службы.....	132
3.3.3	API службы тренировки.....	135
3.3.4	Запуск нового задания на тренировку .....	136
3.3.5	Обновление и доставка информации о статусе задания .....	140
3.3.6	Исходный код тренировки модели классификации намерений .....	142
3.3.7	Управление заданиями на тренировку .....	144
3.3.8	Метрики для устранения неполадок.....	145
3.3.9	Поддержка нового алгоритма или новой версии.....	146
3.4	Тренировочные операторы Kubeflow с открытым исходным кодом .....	147
3.4.1	Тренировочные операторы Kubeflow .....	148
3.4.2	Шаблон оператора/контроллера Kubernetes .....	149
3.4.3	Конструкция тренировочного оператора Kubeflow .....	150
3.4.4	Как использовать тренировочные операторы Kubeflow.....	152
3.4.5	Как интегрировать эти операторы в существующую систему .....	154
3.5	Когда использовать публичное облако.....	155
3.5.1	Когда использовать технологическое решение на базе публичного облака .....	156
3.5.2	Когда создавать свою собственную службу тренировки .....	156
	Резюме .....	158

<b>4</b>	<b>Распределенная тренировка .....</b>	160
4.1	Типы методов распределенной тренировки .....	161
4.2	Параллелизм данных .....	162
4.2.1	Концепция параллелизма данных .....	162
4.2.2	Трудности тренировки с несколькими работниками .....	166
4.2.3	Написание исходного кода распределенной тренировки (с параллелизмом данных) для разных фреймворков тренировки .....	168
4.2.4	Инженерные усилия по распределенной тренировке с параллелизмом данных .....	173
4.3	Образец службы с поддержкой распределенной тренировки с параллелизмом данных .....	176
4.3.1	Обзор службы .....	176
4.3.2	Ознакомление со службой .....	178
4.3.3	Запуск заданий на тренировку .....	180
4.3.4	Обновление и доставка информации о статусе задания ...	184
4.3.5	Конвертация исходного кода тренировки в распределенный режим исполнения .....	185
4.3.6	Улучшения .....	186
4.4	Тренировка больших моделей, не помещающихся на один графический процессор .....	187
4.4.1	Традиционные методы: экономия памяти .....	187
4.4.2	Конвейерный параллелизм модели .....	189
4.4.3	Как разработчикам программного обеспечения поддерживать параллелизм конвейера.....	196
	Резюме .....	197
<b>5</b>	<b>Служба гиперпараметрической оптимизации .....</b>	199
5.1	Понятие гиперпараметров .....	201
5.1.1	Что такое гиперпараметр? .....	201
5.1.2	Причины важности гиперпараметров .....	202
5.2	Понятие гиперпараметрической оптимизации .....	203
5.2.1	Что такое гиперпараметрическая оптимизация? .....	203
5.2.2	Популярные алгоритмы гиперпараметрической оптимизации .....	207
5.2.3	Распространенные подходы к автоматической ГПО.....	214
5.3	Конструирование службы ГПО .....	217
5.3.1	Принципы конструирования службы ГПО .....	217
5.3.2	Общая конструкция службы ГПО .....	219
5.4	Библиотеки ГПО с открытым исходным кодом .....	221
5.4.1	Hyperopt .....	222
5.4.2	Optuna .....	225
5.4.3	Ray Tune .....	227
5.4.4	Следующие шаги .....	231
	Резюме .....	231
<b>6</b>	<b>Конструирование службы раздачи моделей .....</b>	233
6.1	Объяснение процесса раздачи моделей.....	235
6.1.1	Что такое модель машинного обучения? .....	235

6.1.2	Модельное предсказание и модельный вывод .....	237
6.1.3	Что представляет собой раздача модели? .....	238
6.1.4	Трудности раздачи моделей .....	239
6.1.5	Терминология раздачи моделей .....	241
6.2	<b>Распространенные стратегии раздачи моделей .....</b>	242
6.2.1	Прямое встраивание модели .....	242
6.2.2	Служба моделей .....	243
6.2.3	Сервер моделей .....	244
6.3	<b>Конструирование службы предсказания .....</b>	245
6.3.1	Одномодельное приложение .....	246
6.3.2	Многоарендаторное приложение .....	250
6.3.3	Поддержка нескольких приложений в одной системе .....	253
6.3.4	Общие требования к службе предсказания .....	257
	<b>Резюме .....</b>	258
<b>7</b>	<b>Раздача моделей на практике .....</b>	260
7.1	<b>Образец службы моделей .....</b>	261
7.1.1	Ознакомление со службой .....	262
7.1.2	Конструкция службы .....	263
7.1.3	Фронтендовая служба .....	264
7.1.4	Предсказатель классификации намерений .....	271
7.1.5	Вытеснение моделей .....	278
7.2	<b>Образец сервера моделей TorchServe .....</b>	278
7.2.1	Ознакомление со службой .....	279
7.2.2	Конструкция службы .....	280
7.2.3	Фронтендовая служба .....	280
7.2.4	Бэкенд TorchServe .....	281
7.2.5	API TorchServe .....	282
7.2.6	Модельные файлы TorchServe .....	284
7.2.7	Вертикальное масштабирование в Kubernetes .....	289
7.3	<b>Сервер моделей в сопоставлении со службой моделей .....</b>	291
7.4	<b>Экскурсия по инструментам с открытым исходным кодом для раздачи моделей .....</b>	292
7.4.1	TensorFlow Serving .....	293
7.4.2	TorchServe .....	296
7.4.3	Triton Inference Server .....	300
7.4.4	KServe и другие инструменты .....	306
7.4.5	Интеграция инструмента раздачи с существующей системой раздачи .....	307
7.5	<b>Выпуск моделей .....</b>	309
7.5.1	Регистрация моделей .....	311
7.5.2	Загрузка произвольной версии модели в реальном времени с помощью службы предсказания .....	312
7.5.3	Выпуск модели путем обновления дефолтной версии модели .....	314
7.6	<b>Постпроизводственный мониторинг моделей .....</b>	316
7.6.1	Сбор метрических данных и границы качества .....	317
7.6.2	Собираемые метрики .....	317
	<b>Резюме .....</b>	318

<b>8</b>	<b>Склад метаданных и артифактов .....</b>	320
8.1	Введение в артифакты .....	321
8.2	Метаданные в контексте глубокого обучения .....	322
8.2.1	Распространенные категории метаданных .....	323
8.2.2	Зачем нужно управлять метаданными? .....	326
8.3	Конструирование склада метаданных и артифактов .....	329
8.3.1	Принципы конструирования .....	329
8.3.2	Общее конструкционное предложение по складу метаданных и артифактов .....	331
8.4	Технологические решения с открытым исходным кодом .....	334
8.4.1	<i>ML Metadata</i> .....	334
8.4.2	<i>MLflow</i> .....	338
8.4.3	<i>MLflow</i> в сопоставлении с <i>MLMD</i> .....	341
	Резюме .....	342
<b>9</b>	<b>Оркестровка рабочего процесса .....</b>	343
9.1	Введение в оркестровку рабочего процесса .....	344
9.1.1	Что такое рабочий процесс? .....	345
9.1.2	Что такое оркестровка рабочего процесса? .....	346
9.1.3	Трудности использования оркестровки рабочих процессов в глубоком обучении.....	348
9.2	Конструирование системы оркестровки рабочих процессов .....	351
9.2.1	Пользовательские сценарии .....	351
9.2.2	Общая конструкция системы оркестровки .....	354
9.2.3	Принципы конструирования системы оркестровки рабочих процессов.....	356
9.3	Экскурсия по системам оркестровки рабочих процессов с открытым исходным кодом .....	358
9.3.1	<i>Airflow</i> .....	359
9.3.2	<i>Argo Workflows</i> .....	363
9.3.3	<i>Metaflow</i> .....	368
9.3.4	Когда использовать .....	373
	Резюме .....	374
<b>10</b>	<b>Путь к производству .....</b>	375
10.1	Подготовка к продукционализации .....	379
10.1.1	Научные изыскания .....	379
10.1.2	Прототипирование .....	381
10.1.3	Ключевые выводы.....	382
10.2	Продукционализация модели.....	383
10.2.1	Компонентизация исходного кода .....	383
10.2.2	Пакетирование исходного кода.....	385
10.2.3	Регистрация исходного кода.....	386
10.2.4	Настройка рабочего процесса тренировки .....	387
10.2.5	Генерирование модельных выводов .....	388
10.2.6	Интеграция с продуктом .....	389
10.3	Стратегии развертывания моделей .....	390

10.3.1 Канареечное развертывание.....	390
10.3.2 Сине-зеленое развертывание.....	391
10.3.3 Развертывание по принципу работы многорукого бандита .....	392
Резюме .....	393
Дополнение A. Система глубокого обучения «hello world» .....	395
Дополнение B. Экспертиза существующих технологических решений .....	408
Дополнение C. Создание службы гиперпараметрической оптимизации с помощью Kubeflow Katib .....	422
Тематический указатель .....	447

# *Вводное слово*

---

Считается, что система глубокого обучения является эффективной, если она способна наводить мости между двумя разными мирами – научными изысканиями / прототипированием и производственными операциями. Разрабатывающие такие системы коллективы должны уметь общаться с практиками из этих двух миров и работать с разными наборами технических требований и ограничений, которые исходят от каждого из них. В силу этого требуется безупречное понимание конструкционных особенностей компонентов в системах глубокого обучения и принципов их работы в tandemе. Этому аспекту инженерии глубокого обучения посвящено очень мало существующей литературы. Данный информационный пробел становится проблемой, когда технология глубокого обучения внедряется среди младших инженеров программного обеспечения и от них ожидается, что они станут эффективными инженерами в данной сфере.

На протяжении многих лет инженерно-конструкторские коллективы заполняли этот пробел, используя свой приобретенный опыт и извлекая то, что им нужно знать, из литературы. Их работа помогала традиционным инженерам программного обеспечения разрабатывать, конструировать и расширять системы глубокого обучения за относительно короткий промежуток времени. Поэтому я с большим волнением узнал, что Кай и Дональд, оба из которых возглавляли коллективы инженеров глубокого обучения, выступили с очень важной инициативой консолидировать эти знания и поделиться ими в форме книги.

Нам давно пора выпустить всеобъемлющую книгу о разработке систем, которые поддерживают наведение мостов в области глубокого обучения, от научных изысканий и прототипирования до производства. Книга «Разработка систем глубокого обучения», наконец, охватывает эту потребность.

Данная книга начинается с высокоуровневого введения, описывающего суть системы глубокого обучения и ее функциональности.

В последующих главах каждый компонент системы обсуждается подробно, приводится мотивация и дается представление о плюсах и минусах различных конструкционных вариантов.

Каждая глава заканчивается анализом, который помогает читателям оценить варианты, наиболее подходящие для их собственных вариантов использования. В заключение авторы, опираясь на все предыдущие главы, подробно рассказывают о сложном пути перехода от научных изысканий и прототипирования к производству. И для того чтобы помочь инженерам воплотить все эти идеи на практике, они создали образец системы глубокого обучения с полностью рабочим исходным кодом, дабы проиллюстрировать ключевые концепции и предложить попробовать ее тем, кто только начинает работать в этой области.

В целом читатели найдут, что эту книгу легко читать и по ней легко перемещаться, а их понимание способов организации, конструирования и реализации систем глубокого обучения поднимется на совершенно новый уровень. Практики всех уровней знаний, заинтересованные в разработке эффективных систем глубокого обучения, по достоинству оценят эту книгу как бесценный ресурс и справочную информацию. Они прочтут ее один раз, чтобы получить общую картину, а затем будут возвращаться к ней снова и снова при разработке своих собственных систем, конструировании компонентов и принятии важных конструкционных решений, удовлетворяющих все коллективы, которые используют эти системы.

– *Сильвио Саварезе*, вице-президент,

главный научный сотрудник Salesforce

– *Каймин Сюн*, вице-президент Salesforce

# *Предисловие*

---

Чуть более десяти лет назад нам посчастливилось разработать несколько первых продуктовых функциональностей, ориентированных на конечного пользователя, которые были основаны на искусственном интеллекте. Это было грандиозное предприятие. Сбор и организация данных, которые были пригодны для тренировки моделей, в то время не были обычной практикой. Несколько алгоритмов машинного обучения были упакованы в виде готовых к использованию библиотек. Проведение экспериментов требовало ручного управления и разработки конкретно-прикладных рабочих процессов и визуализаций. Для раздачи каждого типа моделей были созданы конкретно-прикладные серверы. За исключением ресурсоемких технологических компаний, почти каждая новая продуктовая функциональность на базе искусственного интеллекта создавалась с нуля. Это была далеко идущая мечта о том, что интеллектуальные приложения однажды станут товаром.

Поработав с несколькими приложениями ИИ, мы поняли, что всякий раз повторяли один и тот же ритуал, и нам показалось, что имеет смысл разработать системный подход с прототипированием для внедрения продуктовых функциональностей на базе ИИ в производство. Результатом этих усилий стал комплект фреймворкового программного обеспечения с открытым исходным кодом под названием PredictionIO, который соединяет в себе самые передовые программные компоненты для сбора и извлечения данных, тренировки и раздачи моделей. Полностью адаптируемый под конкретно-прикладные задачи с помощью API-интерфейсов и пригодный для развертывания в виде служб всего несколькими командами, он помог сократить время, необходимое на каждом этапе – от проведения экспериментов по обработке данных до тренировки и развертывания готовых к производству моделей. Мы были очень рады узнать, что разработчики по всему миру смогли использовать продукт PredictionIO для создания собственных приложений на базе искусственного интеллекта,

что привело к поразительному росту их бизнеса. Позже стартап PredictionIO был приобретен компанией Salesforce для решения аналогичной задачи в еще большем масштабе.

К тому времени, когда мы решили написать эту книгу, индустрия уже процветала благодаря здоровой экосистеме программного обеспечения для искусственного интеллекта. Стало доступно множество алгоритмов и инструментов для решения разных задач. Некоторые поставщики облачных технологий, такие как Amazon, Google и Microsoft, даже предоставляют законченные развернутые системы, которые позволяют коллективам сотрудничать в проведении экспериментов, прототипировании и развертывании производственных приложений в одном централизованном месте. Независимо от вашей цели теперь у вас есть множество вариантов и множество способов их достижения.

Тем не менее, поскольку мы работаем с коллективами над внедрением продуктовых функциональностей глубокого обучения, возникают повторяющиеся вопросы. Почему наша система глубокого обучения сконструирована именно так, а не иначе? Подойдет ли ее конструкция и для других конкретных случаев использования? Мы заметили, что эти вопросы чаще всего задают младшие инженеры программного обеспечения, и мы проинтервьюировали некоторых из них, чтобы выяснить причину. В результате обнаружилось, что их традиционная ученая программа в области разработки программного обеспечения не подготовила их к эффективной работе с системами глубокого обучения. И когда они искали учебные ресурсы, они находили лишь скучную и разрозненную информацию о конкретных системных компонентах, и почти ни в одном ресурсе не обсуждались основы программных компонентов, причины, по которым они были собраны таким образом, и характер их работы вместе, об разуя целостную систему.

В целях решения этой проблемы мы начали разрабатывать базу знаний, которая в конечном итоге превратилась в похожий на руководство учебный материал, объясняющий принципы конструирования каждого компонента системы, плюсы и минусы конструкционных решений, а также обоснование как с технической, так и с продуктовой точки зрения. Мы узнали, что наш материал помогал быстро набирать новых сотрудников в коллективы разработчиков и позволял традиционным инженерам программного обеспечения, не имеющим предшествующего опыта в разработке систем глубокого обучения, входить в курс дела. И решили поделиться этим учебным материалом с гораздо большей аудиторией в формате книги. Мы связались с издательством Manning, и остальное стало историей.

## Отзывы и пожелания

Мы всегда рады отзывам наших читателей. Расскажите нам, что вы думаете об этой книге – что понравилось или, может быть, не по-

нравилось. Отзывы важны для нас, чтобы выпускать книги, которые будут для вас максимально полезны.

Вы можете написать отзыв на нашем сайте [www.dmkpress.com](http://www.dmkpress.com), зайдя на страницу книги и оставив комментарий в разделе «Отзывы и рецензии». Также можно послать письмо главному редактору по адресу [dmkpress@gmail.com](mailto:dmkpress@gmail.com); при этом укажите название книги в теме письма.

Если вы являетесь экспертом в какой-либо области и заинтересованы в написании новой книги, заполните форму на нашем сайте по адресу [http://dmkpress.com/authors/publish\\_book/](http://dmkpress.com/authors/publish_book/) или напишите в издательство по адресу [dmkpress@gmail.com](mailto:dmkpress@gmail.com).

## ***Список опечаток***

Хотя мы приняли все возможные меры для того, чтобы обеспечить высокое качество наших текстов, ошибки все равно случаются. Если вы найдете ошибку в одной из наших книг, мы будем очень благодарны, если вы сообщите о ней главному редактору по адресу [dmkpress@gmail.com](mailto:dmkpress@gmail.com). Сделав это, вы избавите других читателей от недопонимания и поможете нам улучшить последующие издания этой книги.

## ***Нарушение авторских прав***

Пиратство в интернете по-прежнему остается насущной проблемой. Издательства «ДМК Пресс» и Manning Publications очень серьезно относятся к вопросам защиты авторских прав и лицензирования. Если вы столкнетесь в интернете с незаконной публикацией какой-либо из наших книг, пожалуйста, пришлите нам ссылку на интернет-ресурс, чтобы мы могли применить санкции.

Ссылку на подозрительные материалы можно прислать по адресу электронной почты [dmkpress@gmail.com](mailto:dmkpress@gmail.com).

Мы высоко ценим любую помощь по защите наших авторов, благодаря которой мы можем предоставлять вам качественные материалы.

# *Благодарности*

---

Написание книги требует по-настоящему больших персональных усилий, но эта книга была бы невозможна без помощи следующих ниже людей.

В основу значительной части данной книги легла работа с разными коллективами из Salesforce Einstein groups (платформа Einstein, E.ai, Hawking). К этим блестящим и влиятельным коллегам относятся (в алфавитном порядке) Сара Ашер, Джимми Ау, Джон Болл, Аня Бида, Джин Беккер, Ятиш Бхагаван, Джексон Чанг, Химай Десай, Мехмет Эзбидерли, Виталий Гордон, Индира Айер, Арпит Кейл, Шрирам Кришнан, Энни Лэнг, Чан Ли, Эли Левин, Дафна Лю, Лия Макгвайр, Ивайло Михов, Ричард Пак, Генри Сапутра, Рагу Сетти, Шон Сенекал, Карл Скуча, Магнус Торн, Тед Таттл, Иэн Варли, Ян Ян, Марцин Земински и Лео Чжу.

Мы также хотим воспользоваться этой возможностью, чтобы поблагодарить нашего редактора по разработке Фрэнсис Лефковиц. Она не только отличный редактор, который дает потрясающие рекомендации по написанию и прямому редактированию, но и является замечательной наставницей, которая руководила нами на протяжении всего процесса написания книги. Без нее эта книга не была бы такого качества, как сейчас, и не была бы завершена так, как планировалось.

Мы выражаем благодарность коллективу издательства Manning за их руководство на протяжении всего процесса написания книги. Мы по-настоящему ценим возможность узнать мнение читателей на ранних стадиях написания книги благодаря используемой в Manning программе раннего доступа (MEAP).

Всем рецензентам: Алексу Бланку, Амиту Кумару, Аюшу Томару, Бхагвану Коммади, Динкару Джуюлу, Эсрефу Дурне, Гаураву Суду, Гийому Аллеону, Хаммаду Аршаду, Джейми Шафферу, Джапнит Сингх, Джереми Чену, Жуану Динису Феррейре, Кате Паткин, Киту Киму, Ларри Каю, Марии Ане, Микаэлю Дотри, Ник Декруз, Николь

Кенигштейн, Ной Флинн, Оливеру Кортен, Омару Эль Малаку, Пранджалу Ранджану, Рави Суреш Машру, Саиду Эч-Чади, Сандипу Д., Санкету Шарма, Сатеджу Кумару Саху, Саяку Полу, Швете Джоши, Симоне Сгуацца, Шрираму Мачарлу, Сумиту Бхаттачарья, Урсину Стауссу, Видхья Вина и Вэ Луо – ваши предложения помогли сделать эту книгу лучше.

Я хотел бы поблагодарить мою жену Пей Ву за ее безграничную любовь и огромную поддержку на протяжении всего процесса написания этой книги. В трудные времена пандемии Covid Пей оставалась мирным уголком, что позволило писать книгу в окружении оживленной семьи с двумя очаровательными малышами – Кэтрин и Тяньченгом.

Я также хотел бы выразить свою благодарность Ян Сюэ, талантливому разработчику 10X, который написал почти всю лабораторию исходного кода. Его помощь делает ее не только высококачественной, но и простой в освоении. Жена Яна, Донг, всем сердцем его поддерживала, чтобы Ян смог сосредоточиться на лаборатории книги.

Еще один человек, которого я хочу поблагодарить, – Диана Сиболд, талантливый и опытный технический специалист Salesforce. Диана вдохновила меня своим собственным писательским опытом и побудила меня начать писать.

– Кай Ван

Один из основателей стартапа PredictionIO (позже приобретенного компанией Salesforce) научил меня бесценным навыкам разработки продуктов с открытым исходным кодом для разработчиков машинного обучения. Это авантюрное и полезное путешествие было бы невозможным без мужественных душ, которые безмерно доверяли друг другу. Это (в алфавитном порядке) Кеннет Чан, Том Чан, Пэт Феррел, Изабель Ли, Пол Ли, Алекс Мерритт, Томас Стоун, Марко Виверо и Джастин Ип.

Саймон Чан заслуживает особого упоминания. Чан был соучредителем стартапа PredictionIO, и я также имел честь работать с ним и учиться у него в его предыдущих предпринимательских начинаниях. Он был первым человеком, который официально познакомил меня с программированием, когда мы оба учились в одной школе (колледж Ва Янь, Коу-Лун) в Гонконге. Среди других вдохновляющих деятелей школы (в алфавитном порядке) были Дональд Чан, Джейсон Чан, Гамлет Чу, Ка Куен Фу, Джеффри Хау, Фрэнсис Конг, Эрик Лау, Кам Лау, Рэйлекс Ли, Кевин Лей, Дэнни Шинг, Норман Со, Стивен Танг и Лобо Вонг.

Я чрезвычайно благодарен своим родителям и моему брату Рональду. Они помогли мне познакомиться с компьютерами в раннем возрасте. Их постоянная поддержка сыграла жизненно важную роль в годы моего становления, когда я стремился стать инженером-компьютерщиком.

Мой сын Спенсер является ходячим доказательством того, почему биологические глубокие нейронные сети – самые удивительные вещи в мире. Он – чудесный дар, который каждый день показывает мне, что я всегда могу расти и становиться лучше.

Словами не передать, как много значит для меня моя жена Вики. Она всегда способна пробудить во мне все самое лучшее, чтобы в трудные моменты я мог продолжать двигаться вперед. Она – лучшая компаньонка, о которой я когда-либо мог мечтать.

– *Дональд Сзето*

# *О книге*

---

Данная книга призвана оснастить инженеров знаниями о том, как конструировать, строить или организовывать эффективные системы машинного обучения и адаптировать эти системы под любые потребности и ситуации, с которыми они могут столкнуться. Разрабатываемые ими системы будут облегчать, автоматизировать и ускорять разработку проектов машинного обучения (в частности, проектов глубокого обучения) в различных областях.

В области глубокого обучения именно модели привлекают все внимание. Возможно, это справедливо, если учесть, что на рынок регулярно поступают новые приложения, разработанные на основе этих моделей, – приложения, приводящие потребителей в восторг, такие как камеры слежения, распознающие человека, виртуальные персонажи в интернет-videоиграх, которые ведут себя как настоящие люди, программа, которая может писать исходный код решения поставленных перед ней произвольных задач, а также передовые системы помощи водителю, которые в один прекрасный день приведут к полностью автономным и самоуправляемым автомобилям. За очень короткий промежуток времени область глубокого обучения наполнилась огромным волнением и многообещающим потенциалом, ожидающим полного воплощения.

Но модель действует не в одиночку. Материализация продукта или службы предусматривает, что модель будет располагаться в системе или на платформе (мы используем эти термины взаимозаменяемо), которая поддерживает модель различными службами и хранилищами. Для этого, среди прочего, нужны, например, API, менеджер наборов данных и склад артифактов<sup>1</sup> и метаданных. Таким образом, за каждым коллективом разработчиков моделей глубокого обучения

---

<sup>1</sup> Артифакт – это продукт деятельности вообще, в данном случае результат работы системы или модели, тогда как артефакт – это произведение искусства или изделие художественного творчества. – Прим. перев.

стоит коллектив разработчиков, в чьи обязанности входит не само глубокое обучение, а создание инфраструктуры, которая содержит модель и все остальные компоненты.

Мы заметили в отрасли проблему, которая заключается в том, что нередко разработчики, которым поручено разрабатывать систему и компоненты глубокого обучения, обладают лишь поверхностными знаниями о глубоком обучении. Они не понимают конкретных требований, предъявляемых к инженерно-конструкторской работе над системой в глубоком обучении, поэтому при разработке системы они склонны следовать обобщенным подходам. Например, они могут принять решение передать всю работу, связанную с разработкой модели глубокого обучения, исследователю данных и сосредоточиться только на автоматизации. И поэтому разрабатываемая ими система будет опираться на традиционную систему планирования списков заданий или бизнес-аналитическую систему анализа данных, которая не оптимизирована ни под характер выполнения заданий глубокого обучения, ни под специфичные для глубокого обучения шаблоны доступа к данным. В результате систему будет сложно использовать для разработки моделей, а скорость доставки моделей останется низкой. По сути, инженеров, которым не хватает глубокого понимания глубокого обучения, просят разрабатывать системы с поддержкой моделей глубокого обучения. Как следствие эти инженерные системы неэффективны и не подходят для систем глубокого обучения.

О разработке моделей глубокого обучения было написано немало с точки зрения исследователя данных, охватывая сбор данных и насыщение наборов данных, написание алгоритмов тренировки и т. п. Но очень немногие книги или даже блоги были посвящены системе и службам, которые поддерживают все эти мероприятия по глубокому обучению.

В данной книге мы обсуждаем строительство и конструирование систем глубокого обучения с точки зрения разработчика программного обеспечения. Наш подход заключается в том, чтобы сначала описать типичную систему глубокого обучения в целом, включая ее главнейшие компоненты и их взаимосвязь; затем в отдельной главе мы углубляемся в каждый главнейший компонент. Мы начинаем каждую главу о компонентах с изложения технических требований. Затем знакомим с принципами конструирования, а также образцами служб / исходного кода и, наконец, оцениваем технологические решения с открытым исходным кодом.

Поскольку мы не можем охватить в книге все существующие системы глубокого обучения (от производителя или с открытым исходным кодом), то сосредоточиваемся на изложении технических требований и принципов конструирования (с примерами). Изучив эти принципы, попробовав примеры служб из книги и прочитав наше обсуждение вариантов с открытым исходным кодом, мы надеемся, что читатели смогут провести собственное исследование и найти то, которое подходит им лучше всего.

## *Кому следует прочитать эту книгу?*

Первостепенная аудитория этой книги – инженеры программного обеспечения (включая недавно окончивших учебу студентов со специализацией в области вычислительных наук), которые хотят быстро перейти к разработке систем глубокого обучения, например те, кто хочет работать на платформах глубокого обучения или интегрировать какую-либо функциональность искусственного интеллекта – например, раздачу моделей – в свои продукты.

Исследователи данных, инженеры-изыскатели, менеджеры и все остальные, кто использует машинное обучение для решения реальных задач, также найдут эту книгу полезной. Разобравшись в опорной инфраструктуре (или системе), они будут оснащены всем необходимым для предоставления точной обратной связи коллективу инженеров в том, что касается повышения эффективности процесса разработки моделей.

Это инженерная книга, и вам не нужны знания в области машинного обучения, но вы должны быть знакомы с базовыми концепциями вычислительных наук и инструментами программирования, такими как микросервисы, gRPC и Docker, чтобы работать с лабораторией исходного кода и понимать технический материал. Независимо от вашего образования, вы все равно сможете извлечь пользу из нетехнических материалов книги, которые помогут вам лучше понять принципы работы машинного обучения и систем глубокого обучения, для того чтобы переносить продукты и службы из области идей в производство.

Прочитав эту книгу, вы сможете понять механизмы работы систем глубокого обучения и способы разработки каждого компонента. Вы также поймете ситуации, когда следует собирать технические требования пользователей, транслировать их в конструкционные варианты системных компонентов и интегрировать компоненты с целью создания целостной системы, которая поможет вашим пользователям быстро разрабатывать и предоставлять функциональные возможности глубокого обучения.

## *Как эта книга организована: дорожная карта*

В этой книге 10 глав и три дополнения (включая одно лабораторное дополнение). В первой главе объясняется, что такое цикл освоения глубокого обучения и как выглядит базовая система глубокого обучения. В следующих главах мы подробно рассмотрим каждый функциональный компонент эталонной системы глубокого обучения. Наконец, в последней главе обсуждается вопрос о том, как модели отправляются в производство. Дополнение к книге содержит лабораторию исходного кода, позволяющую читателям опробовать образец системы глубокого обучения.

В главе 1 описывается, что такое система глубокого обучения, разные участвующие в разработке системы интересанты и как они

взаимодействуют с ней в целях ее оснащения функциональными возможностями глубокого обучения. Это взаимодействие называется циклом освоения глубокого обучения. В добавок вы сформируете для себя концепцию системы глубокого обучения, именуемую эталонной системной архитектурой, которая содержит все необходимые элементы и может быть адаптирована в соответствии с вашими требованиями.

Главы со 2 по 9 охватывают каждый стержневой компонент эталонной архитектуры системы глубокого обучения, такой как служба управления наборами данных, служба тренировки моделей, служба автоматической гиперпараметрической оптимизации и служба оркестровки рабочих процессов.

В главе 10 обсуждается вопрос о том, как подготовить конечный продукт в фазе научного изыскания или прототипирования к выпуску в публичное пространство. В дополнении А представлен образец системы глубокого обучения и демонстрируется лабораторное упражнение, в дополнении В приводится обзор существующих технологических решений, а в дополнении С обсуждается система Kubeflow Katib.

## Об исходном коде

Мы считаем, что самый лучший способ учиться – это делать, практиковаться и экспериментировать. В целях демонстрации описанных в этой книге принципов конструирования и получения практического опыта мы создали образец системы глубокого обучения и лабораторию исходного кода. Весь исходный код, инструкции по настройке и лабораторные скрипты образца системы глубокого обучения доступны на GitHub (<https://github.com/orca3/MiniAutoML>). Вы также можете получить исполняемые фрагменты исходного кода из онлайн-версии этой книги в liveBook по адресу <https://livebook.manning.com/book/software-engineers-guide-to-deep-learning-system-design> и с веб-сайта издательства Manning ([www.manning.com](http://www.manning.com)).

Лаборатория исходного кода «hello world» (в дополнении А) содержит полную, хотя и упрощенную, мини-систему глубокого обучения с наиболее важными компонентами (управлением наборами данных, тренировкой и раздачей моделей). Мы предлагаем вам опробовать указанную там систему после прочтения первой главы книги либо сделать это до того, как вы испытаете наши образцы служб, описанные в этой книге. Данная лаборатория также предоставляет скрипты командной оболочки и ссылки на все ресурсы, необходимые для того, чтобы приступить к работе.

Помимо лаборатории исходного кода, эта книга содержит множество примеров исходного кода в виде отдельных нумерованных листингов и внутри обычного текста. В обоих случаях исходный код отформатирован шрифтом фиксированной ширины, подобным этому, чтобы отделять его от обычного текста. Иногда исходный код также выделяется **жирным шрифтом**, чтобы выделять тот исходный код, который

изменился по сравнению с предыдущими шагами в данной главе, например когда новая функциональная возможность добавляется к существующей строке исходного кода.

Во многих случаях изначальный исходный код был переформатирован; мы добавили переносы строк и переработали отступы, чтобы уместиться в доступное пространство страницы книги. В редких случаях даже этого было недостаточно, и листинги включали маркеры продолжения строки (➡). Вдобавок нередко комментарии в исходном коде из листингов удалялись, когда исходный код описывался в тексте. Многие листинги сопровождаются аннотациями к исходному коду, выделяющими важные концепции.

Конец ознакомительного фрагмента.  
Приобрести книгу можно  
в интернет-магазине  
«Электронный универс»  
[e-Univers.ru](http://e-Univers.ru)