

ШИГЛА

ЦИФРОВАЯ
ИНДУСТРИЯ
ПРОМЫШЛЕННОЙ
РОССИИ

ЦИПР

6-8 июня
Иннополис

2018

О ПРОБЛЕМАХ ИНТЕЛЛЕКТУАЛЬНОГО АНАЛИЗА БОЛЬШИХ ДАННЫХ

*Чехович Юрий Викторович, зав. отделом
Интеллектуальных систем ФИЦ ИУ РАН, к.ф.-м.н.*

НЕСКОЛЬКО СЛОВ О НАШЕЙ КОМАНДЕ

- Вычислительный центр РАН
- Кафедры МФТИ и ВМК МГУ
- Несколько внедряющих компаний:
Форексис, Гудфоркаст, Антиплагиат
- Большой опыт решения задач анализа данных



КЛЮЧЕВОЙ ФАКТОР В АНАЛИЗЕ БОЛЬШИХ ДАННЫХ

Что главное

Для успешного создания и развития технологий хранения и анализа больших данных необходимо выделить ключевые факторы будущего успеха:

- Математические методы
- Оборудование для хранения, обработки и передачи
- Подготовка специалистов
- Кооперация



КЛЮЧЕВОЙ ФАКТОР В АНАЛИЗЕ БОЛЬШИХ ДАННЫХ

Что главное

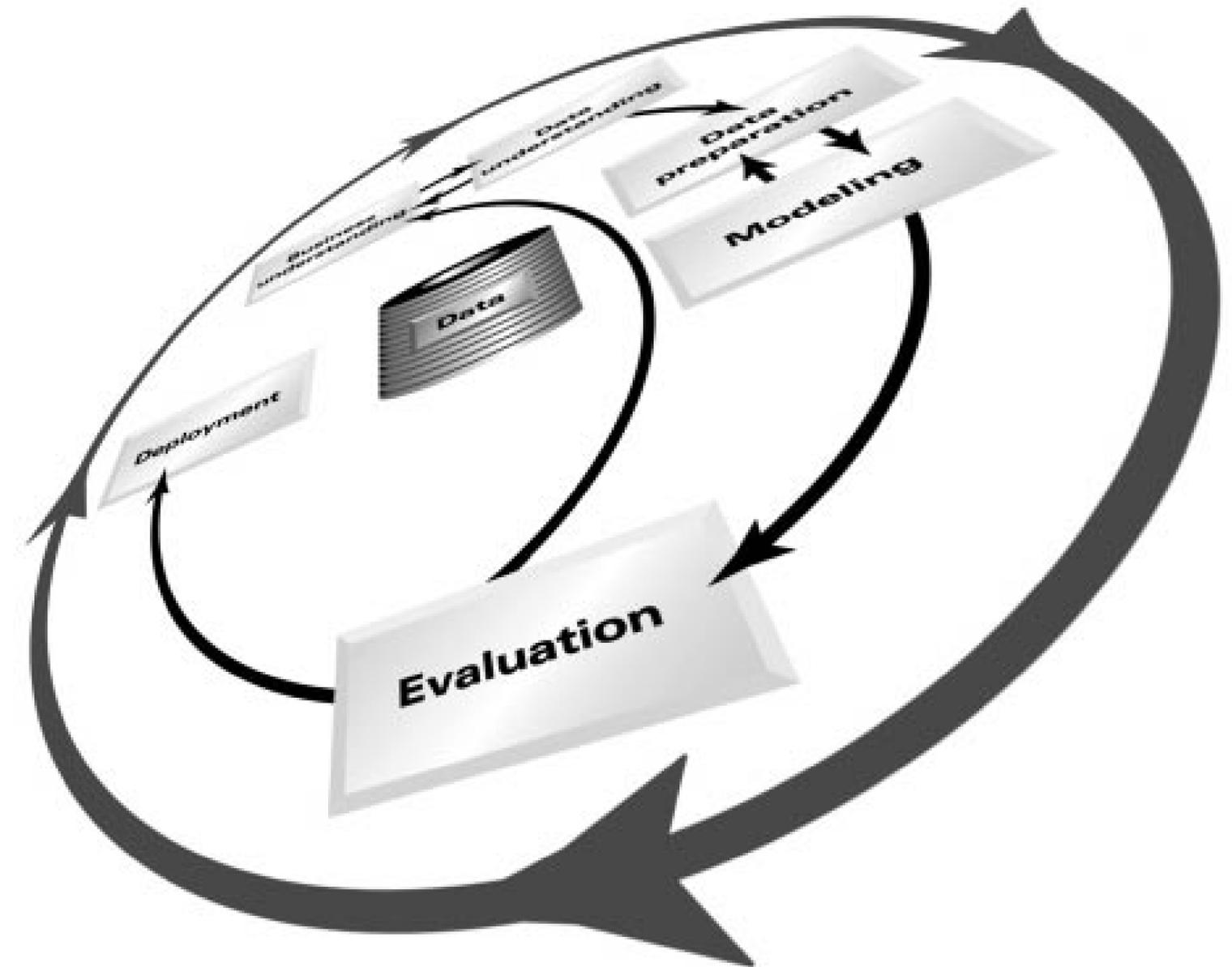
Для успешного создания и развития технологий хранения и анализа больших данных необходимо выделить ключевые факторы будущего успеха:

- Математические методы
- Оборудование для хранения, обработки и передачи
- Подготовка специалистов
- Кооперация

Методология!

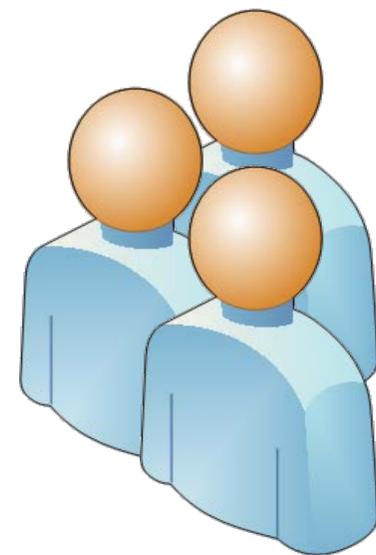
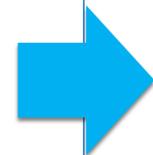
CRISP-DM

Методология ведения промышленных проектов интеллектуального анализа данных



ИСТОРИЯ И ОСОБЕННОСТИ CRISP-DM

Teradata
a division of NCR



SIG (Special Interest Group)



1996
Начали
задумываться

1997-1999
Проведение конференций,
анализ возможных сценариев

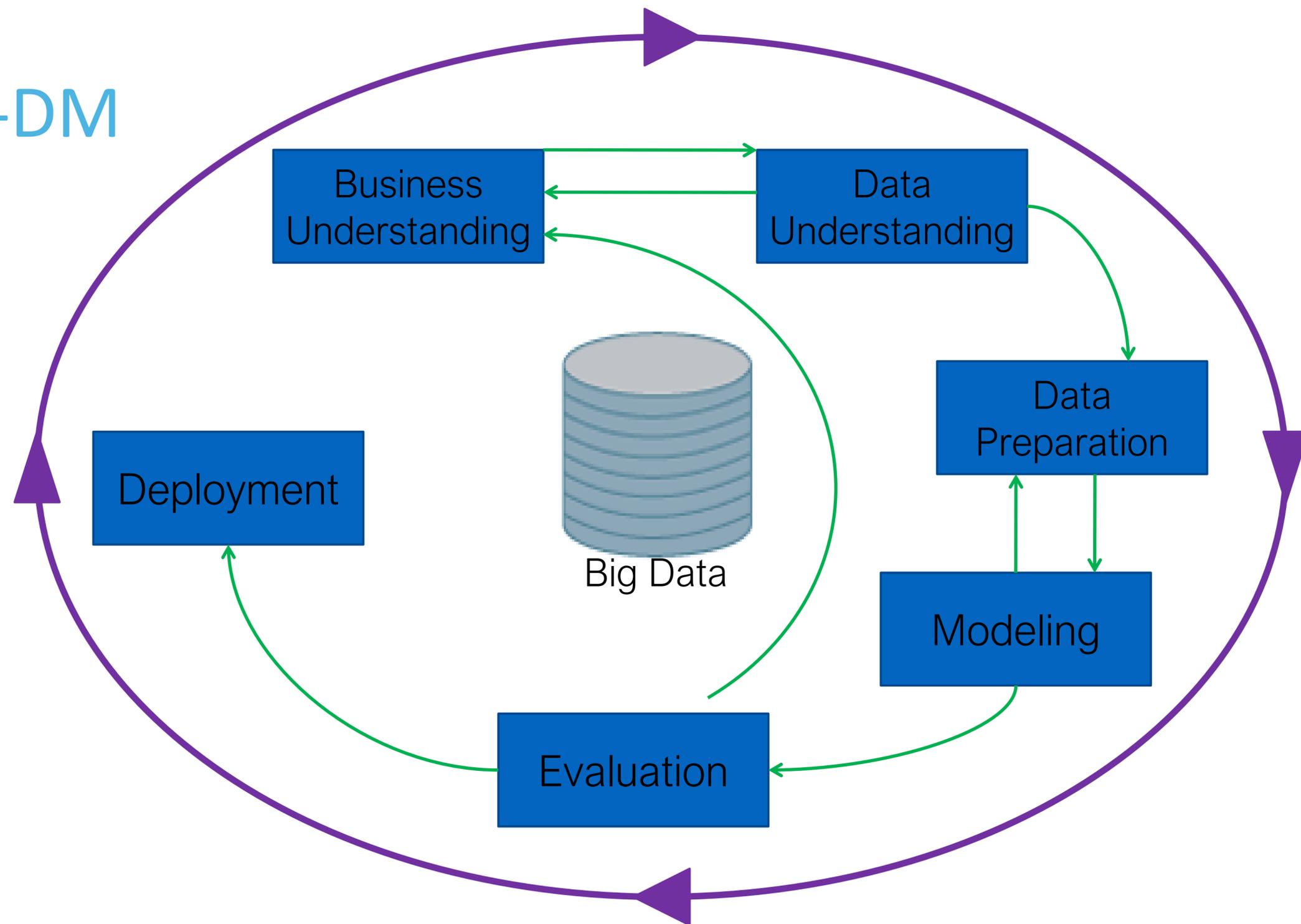
2000
CRISP-DM 1.0

ИСТОРИЯ И ОСОБЕННОСТИ CRISP-DM

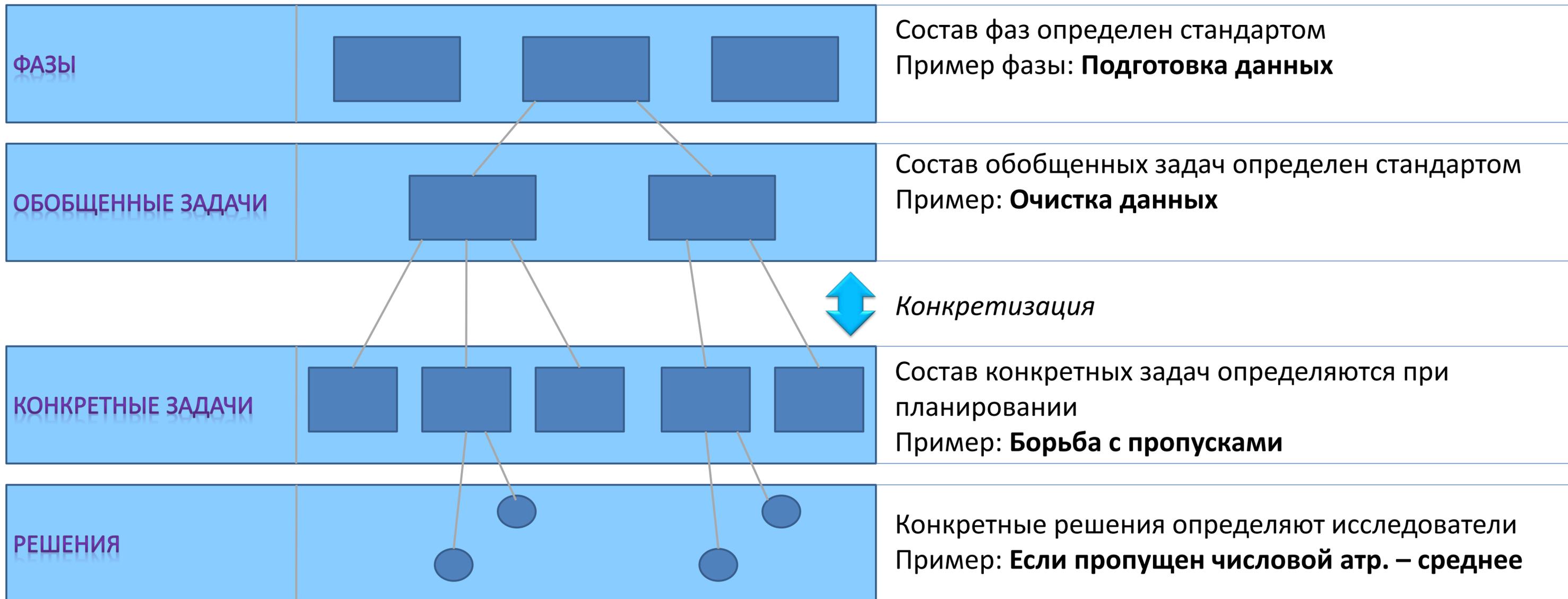
- ④ Стандартный процесс для ведения исследовательских проектов в области анализа данных
- ④ Не зависит от прикладной области
- ④ Не зависит от используемого ПО
- ④ Не зависит от применяемых алгоритмов и решаемых задач

ОСНОВНЫЕ ПРИНЦИПЫ CRISP-DM

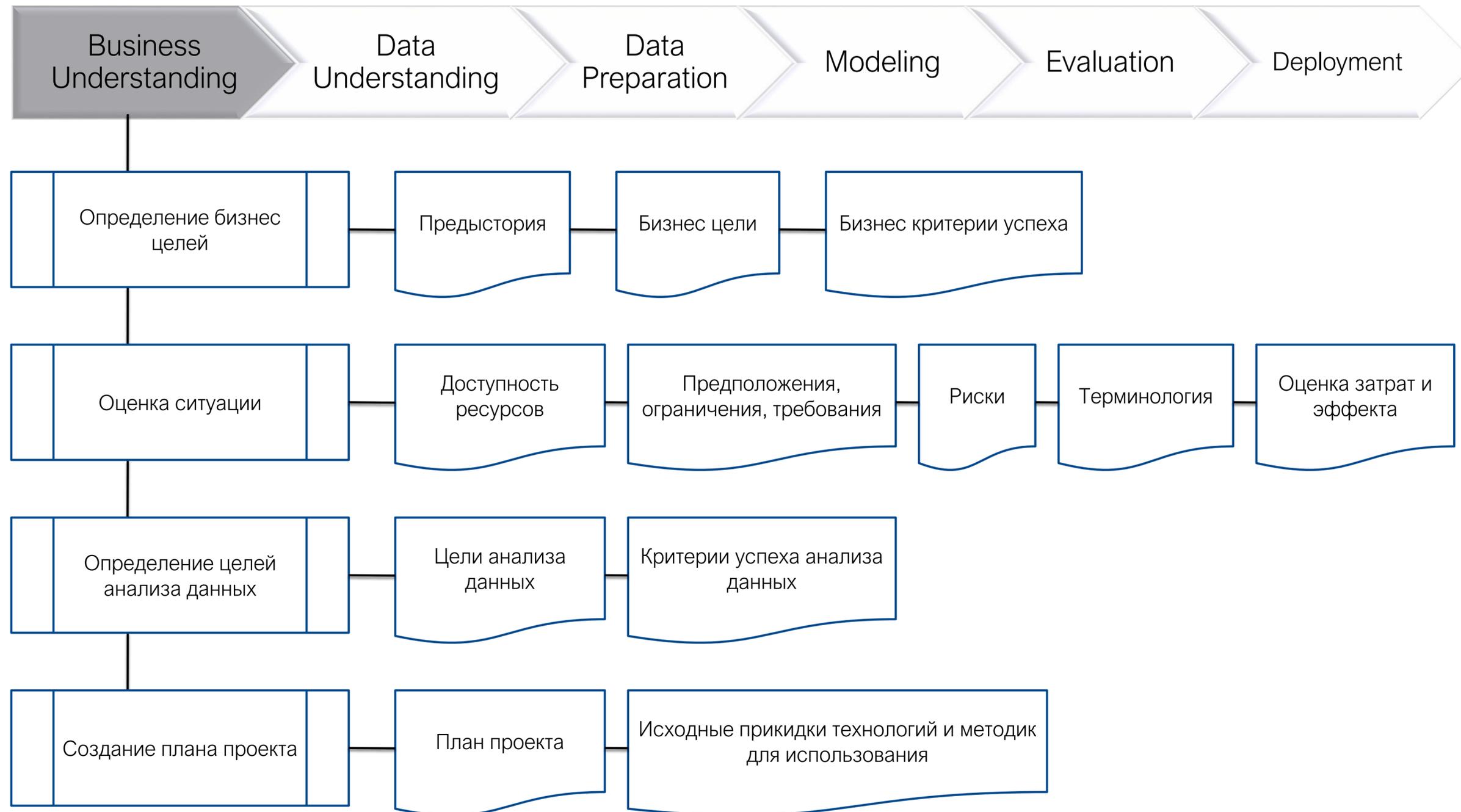
Фазы CRISP-DM



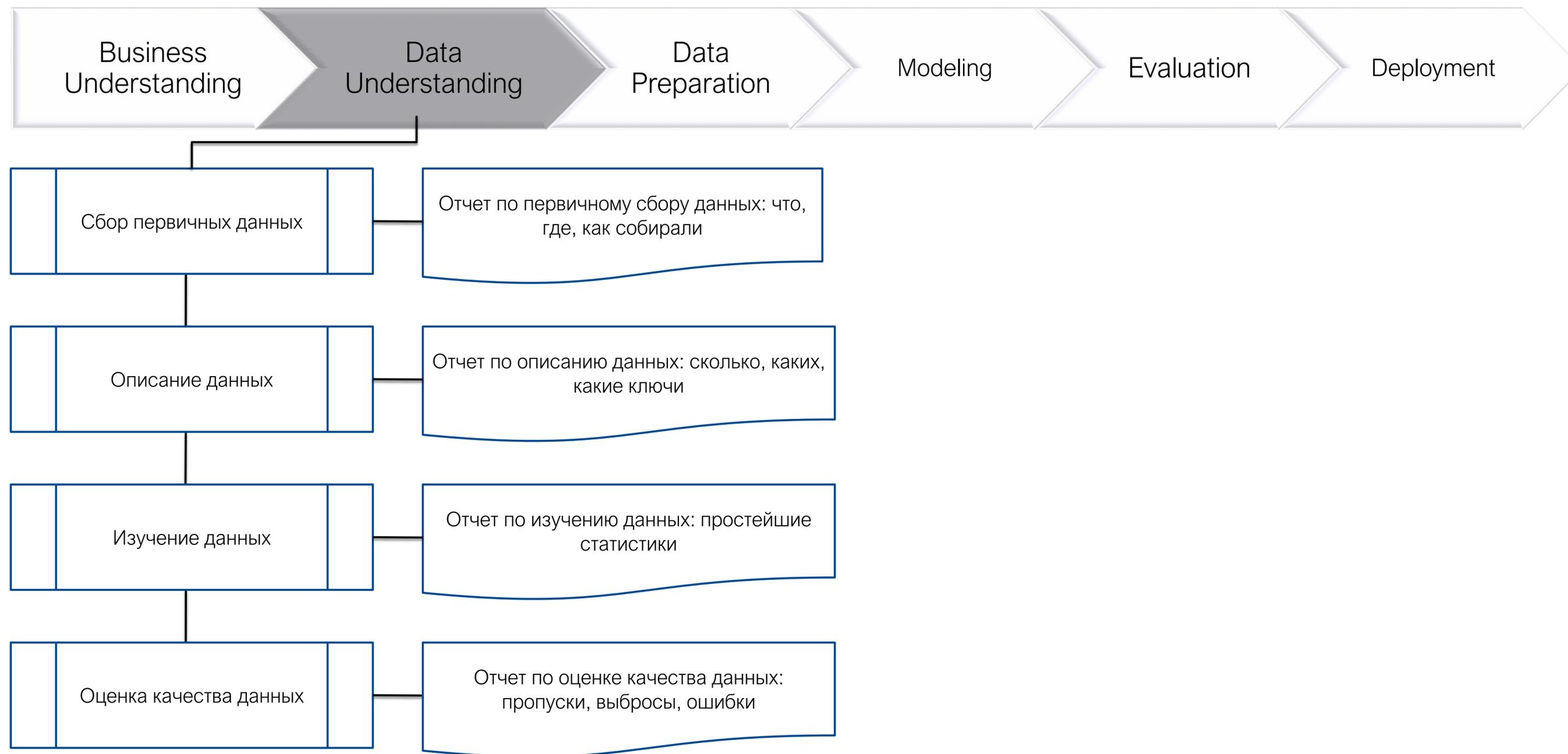
Иерархическая декомпозиция



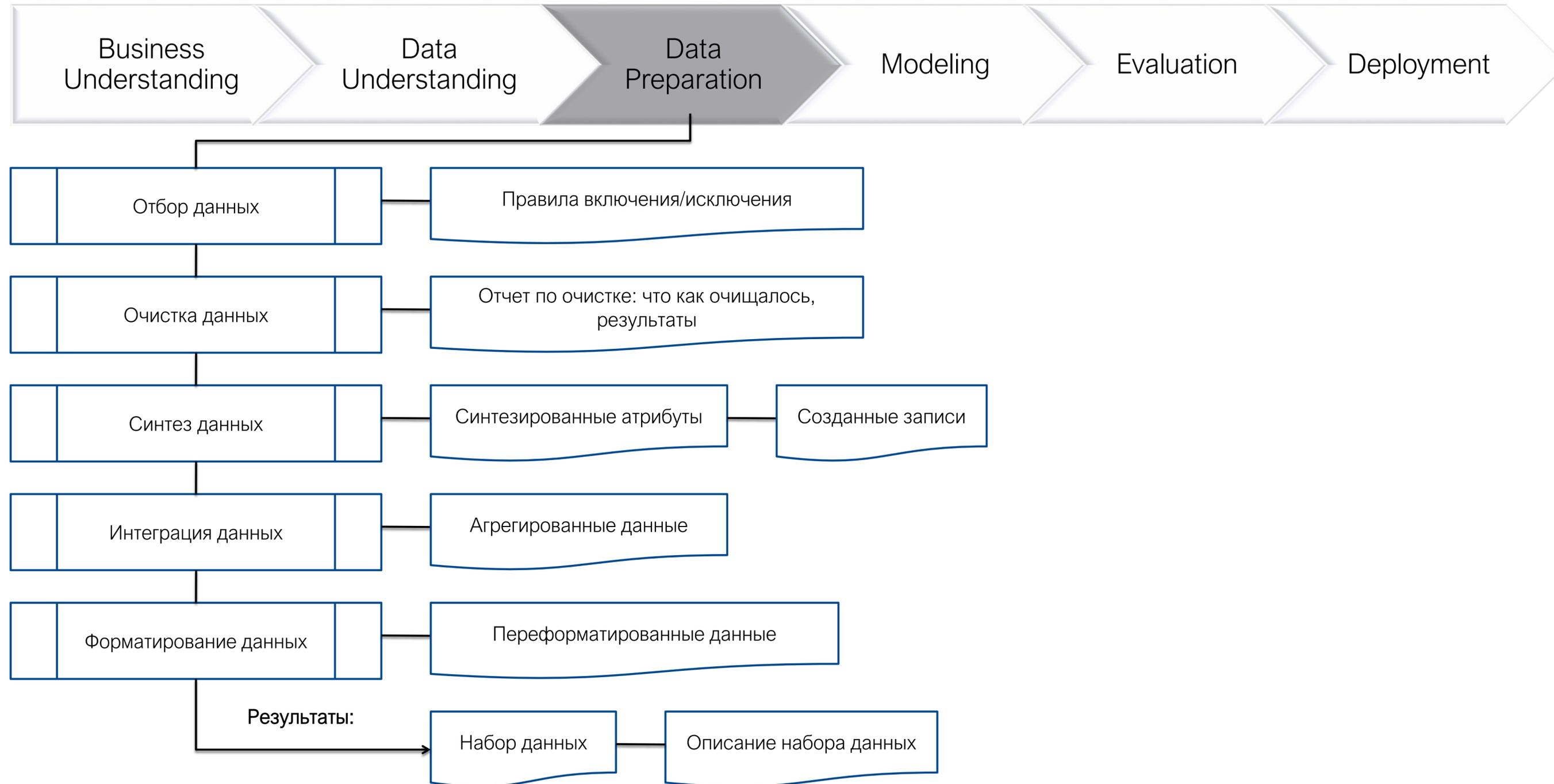
СОДЕРЖАТЕЛЬНАЯ ПОСТАНОВКА ЗАДАЧИ (BUSINESS UNDERSTANDING)



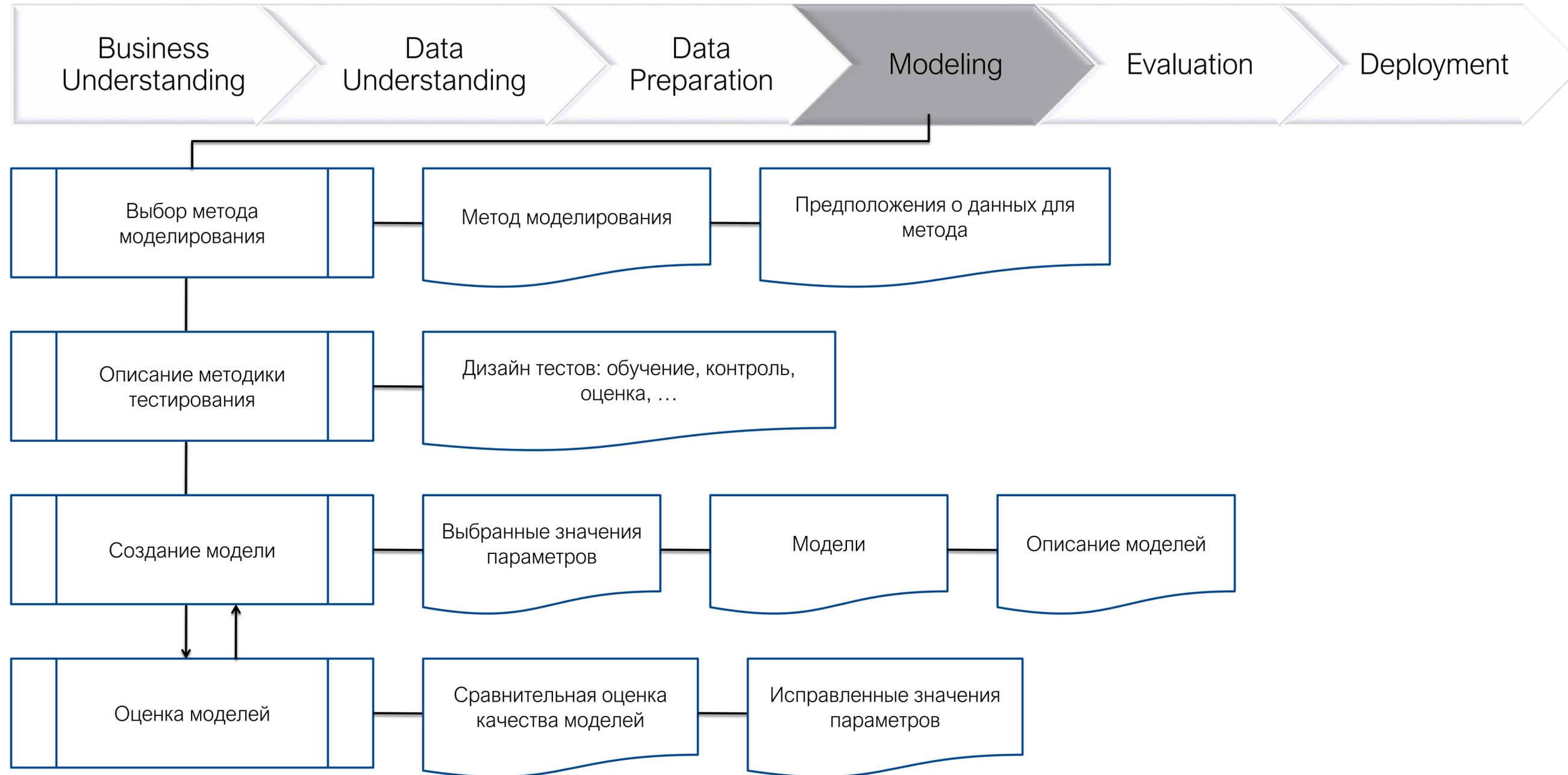
АНАЛИЗ ИСТОЧНИКОВ ДАННЫХ (DATA UNDERSTANDING)



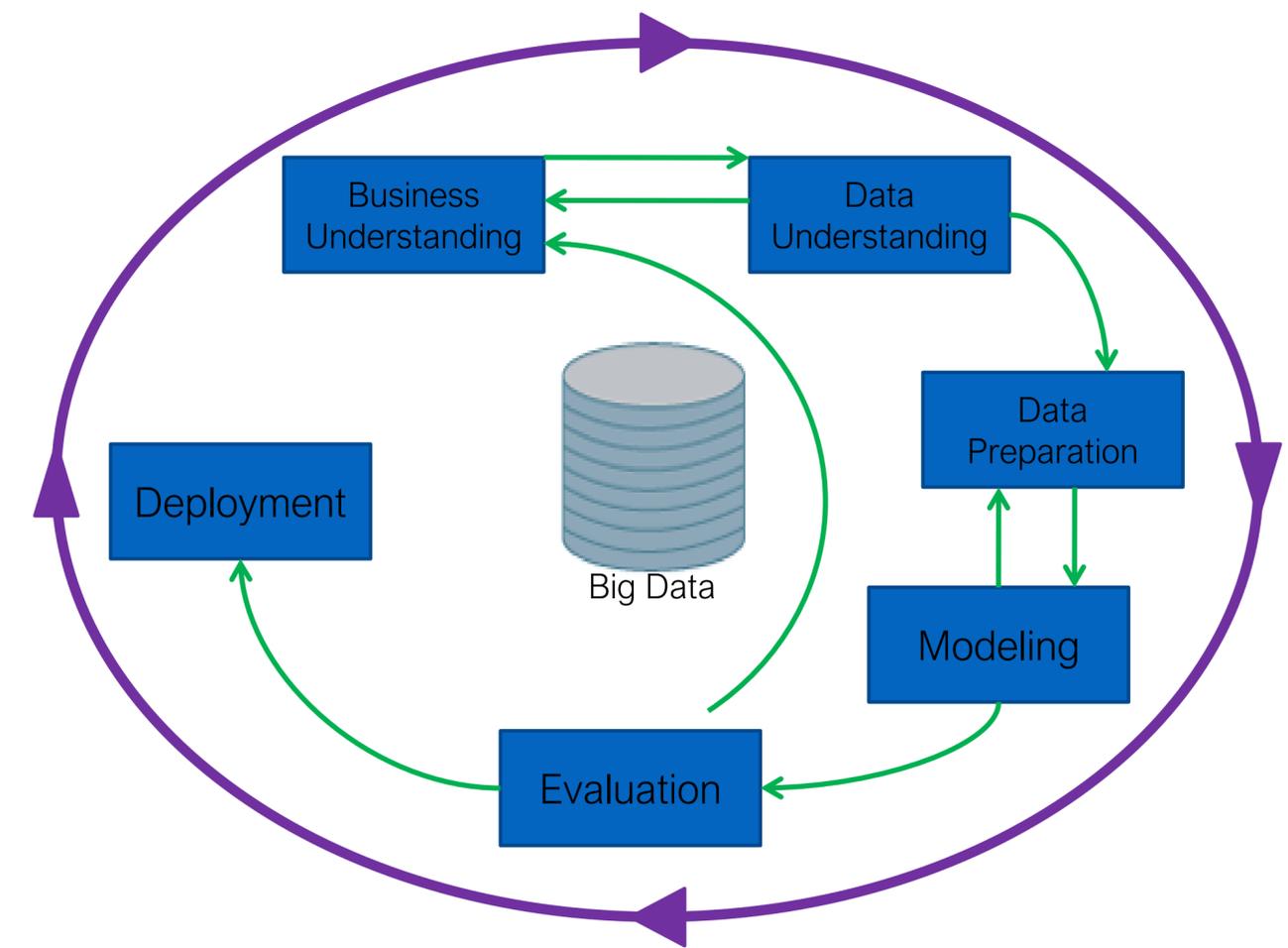
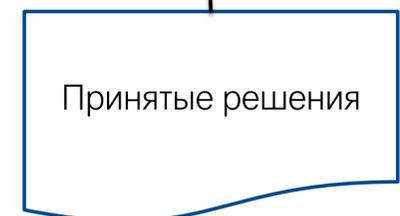
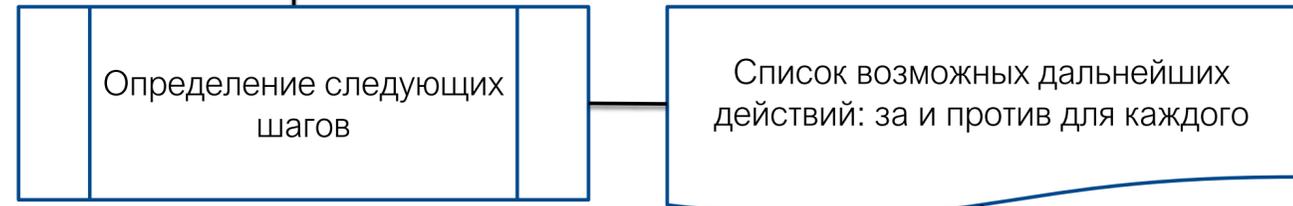
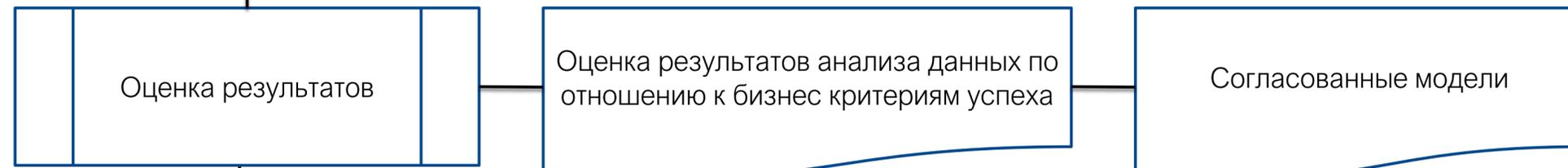
ПРЕДОБРАБОТКА ДАННЫХ (DATA PREPARATION)



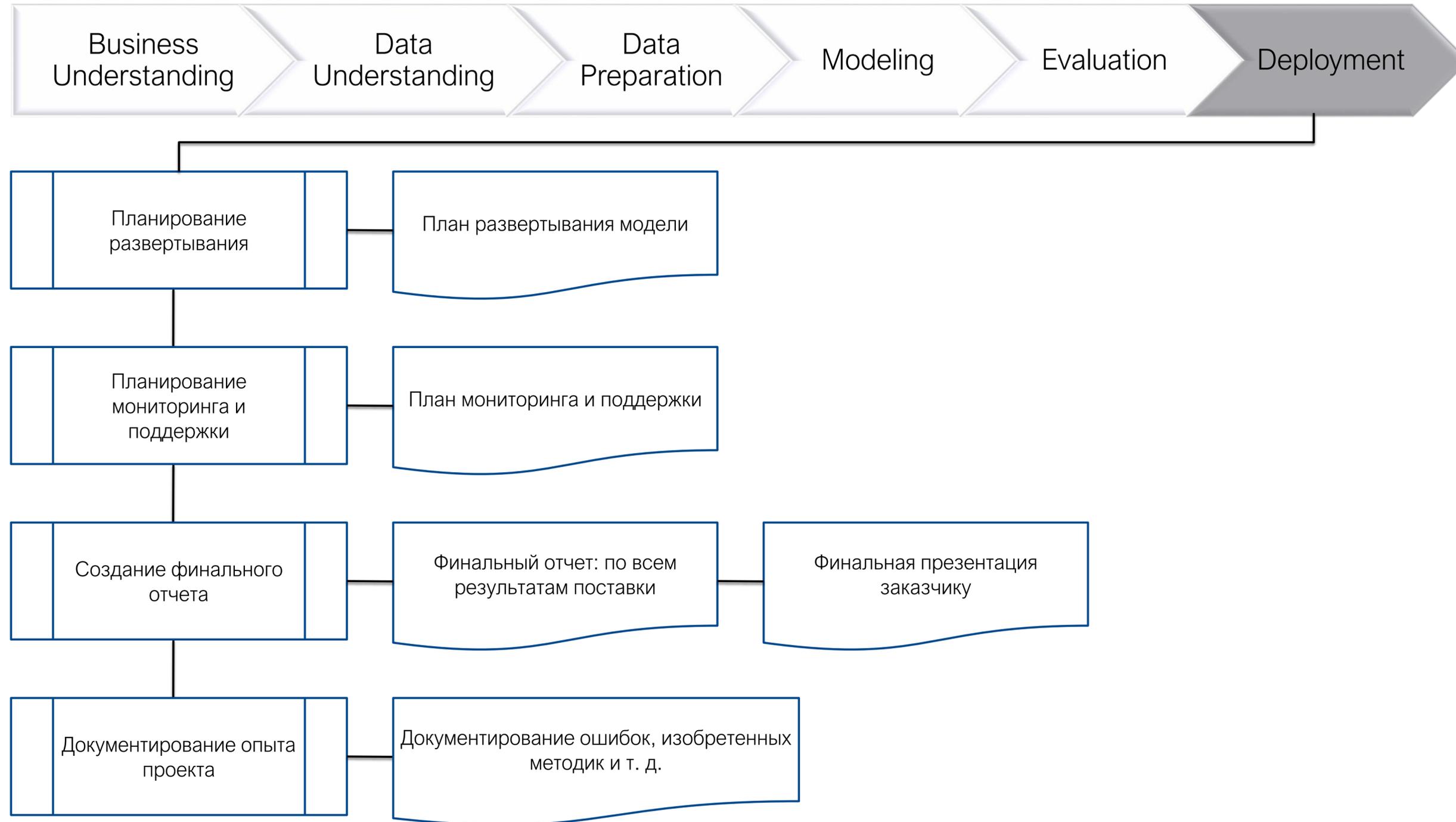
ПОИСК РЕШЕНИЯ (MODELING)



ОЦЕНКА РЕШЕНИЯ (EVALUATION)



РАЗВЕРТЫВАНИЕ (DEPLOYMENT)



ЗАМЕЧАНИЯ И ВЫВОДЫ

- ⦿ CRISP-DM не панацея
- ⦿ CRISP-DM не мешает таланту
- ⦿ CRISP-DM позволяет накапливать и передавать опыт
- ⦿ CRISP-DM создавался как стандарт анализа данных. Большие данные не проблема?
- ⦿ Не проблема!

Спасибо за внимание!