

---

# Сотрудничество ЛНСб с Яндекс

---

ПОДГОТОВИЛ:  
БИКБАЕВ Т.Э.

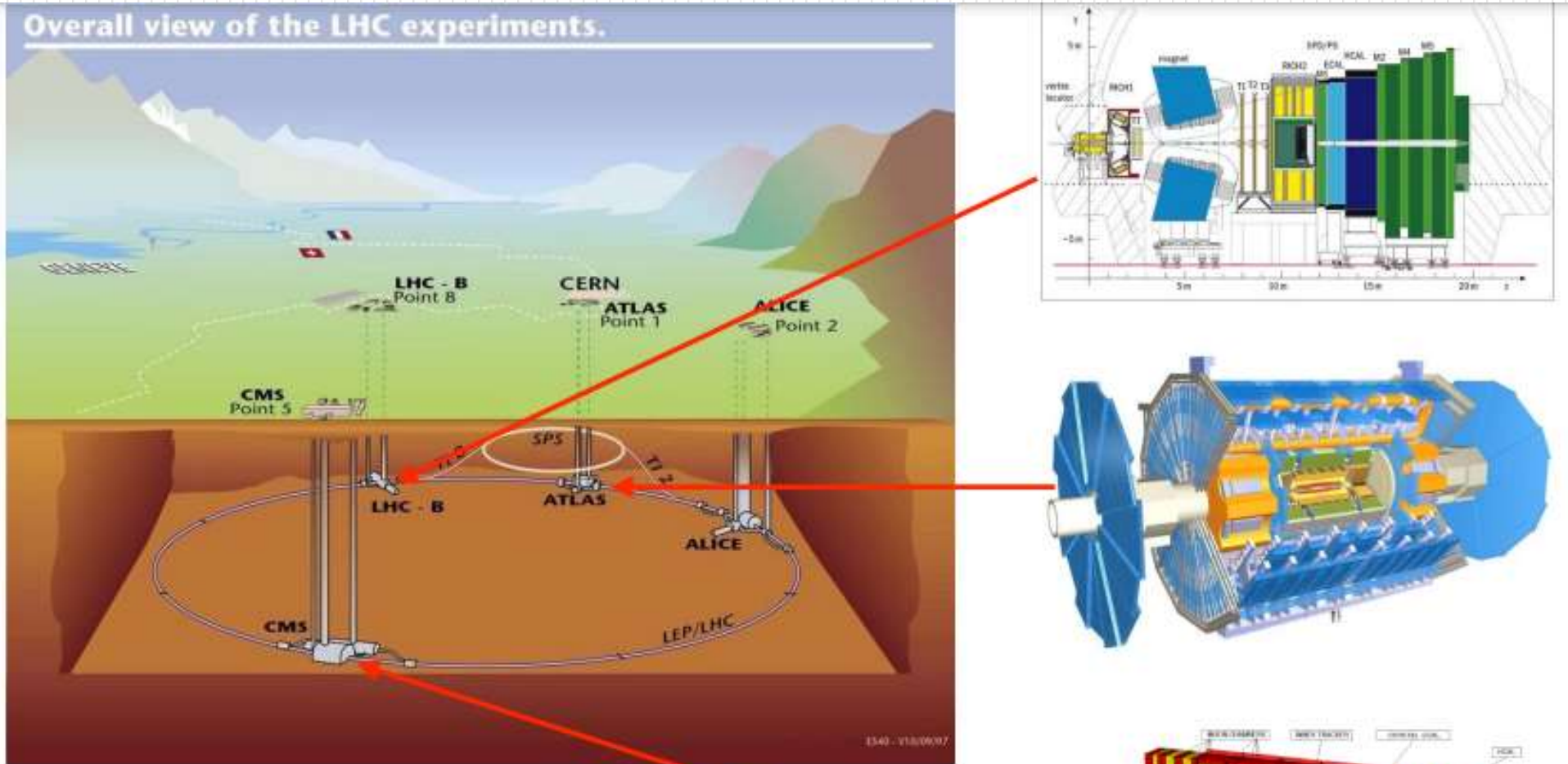
**Яндекс**

# Что такое Яндекс?

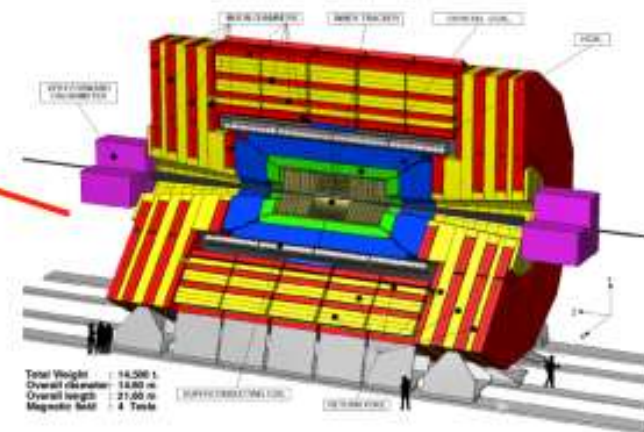
 <b>Поиск</b> Ответы на любые вопросы	 <b>Картинки</b> Изображения всех цветов и размеров	 <b>Видео</b> Просмотр фильмов, сериалов, телешоу, музыкальных роликов
 <b>Новости</b> Картина дня, созданная автоматически	 <b>Погода</b> Прогноз в вашем городе и по всему миру	 <b>Карты</b> Подробные схемы городов, маршруты без пробок
 <b>Почта</b> Электронный ящик без спама и вирусов	 <b>Маркет</b> Товары, сравнение цен, отзывы покупателей	 <b>Яндекс.Браузер</b> Простой и безопасный интернет
 <b>Афиша</b> Развлекательные мероприятия	 <b>Такси</b> Свободные водители поблизости	 <b>Музыка</b> Персональные рекомендации
 <b>Деньги</b> Онлайн-платежи и электронный кошелёк	 <b>Диск</b> Безопасное облако для ваших файлов	 <b>Недвижимость</b> Объявления о комнатах, квартирах и домах
 <b>Авто.ру</b> Огромный выбор новых и подержанных автомобилей	 <b>Авиабилеты</b> Большой выбор предложений от авиакомпаний и агентств	 <b>Работа</b> Подбор вакансий с популярных сайтов поиска работы

- Технологии поиска
- Технологии машинного обучения
- Технологии работы с большими объемами данных

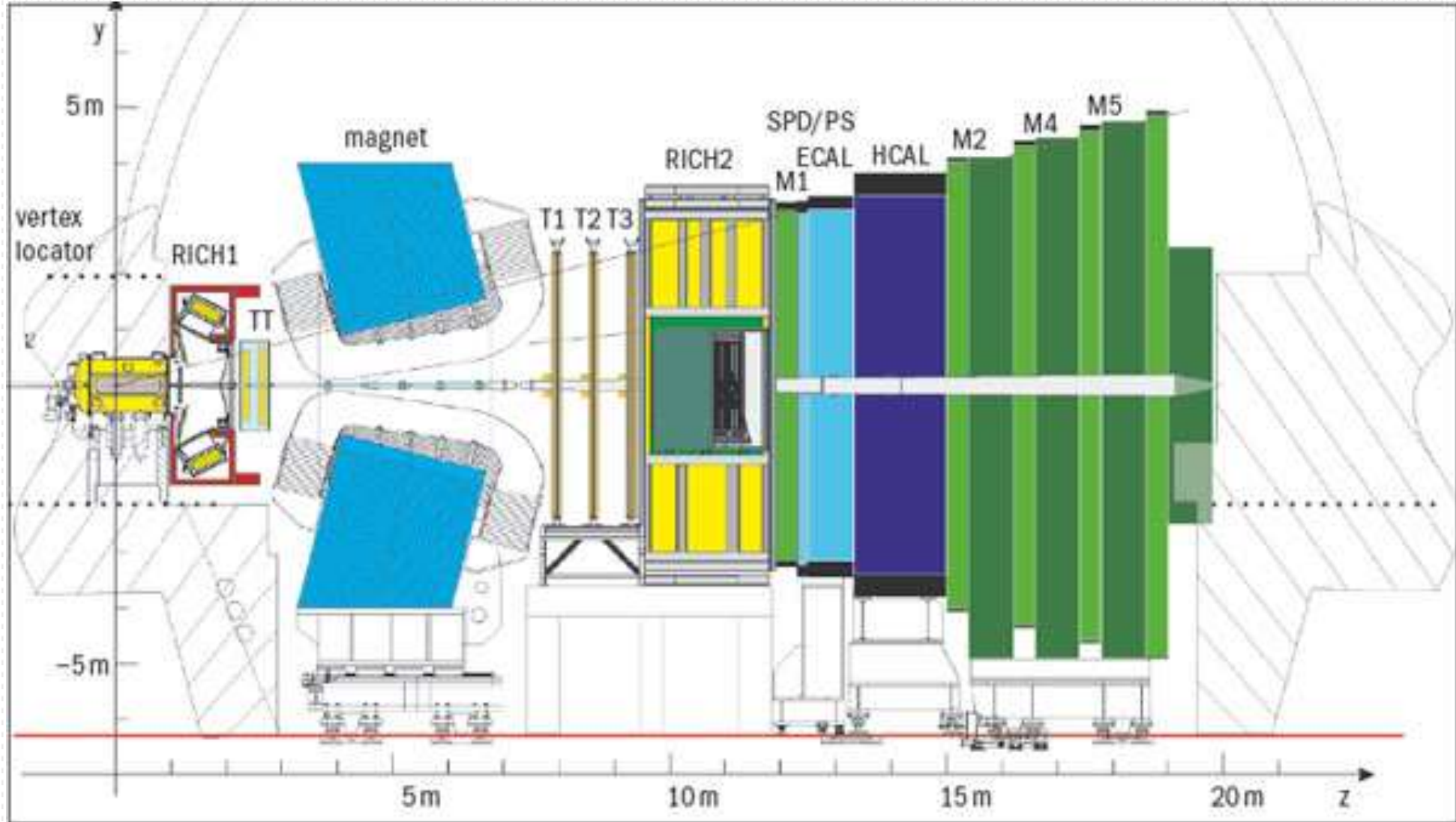
# Большой Адронный Коллайдер (LHC)



- ◇ 2 протонных пучка, скорость - 99.99999991% скорости света
- ◇ 4 точки пересечения пучков
- ◇ энергия в пучках ~ энергии Титаника на полном ходу

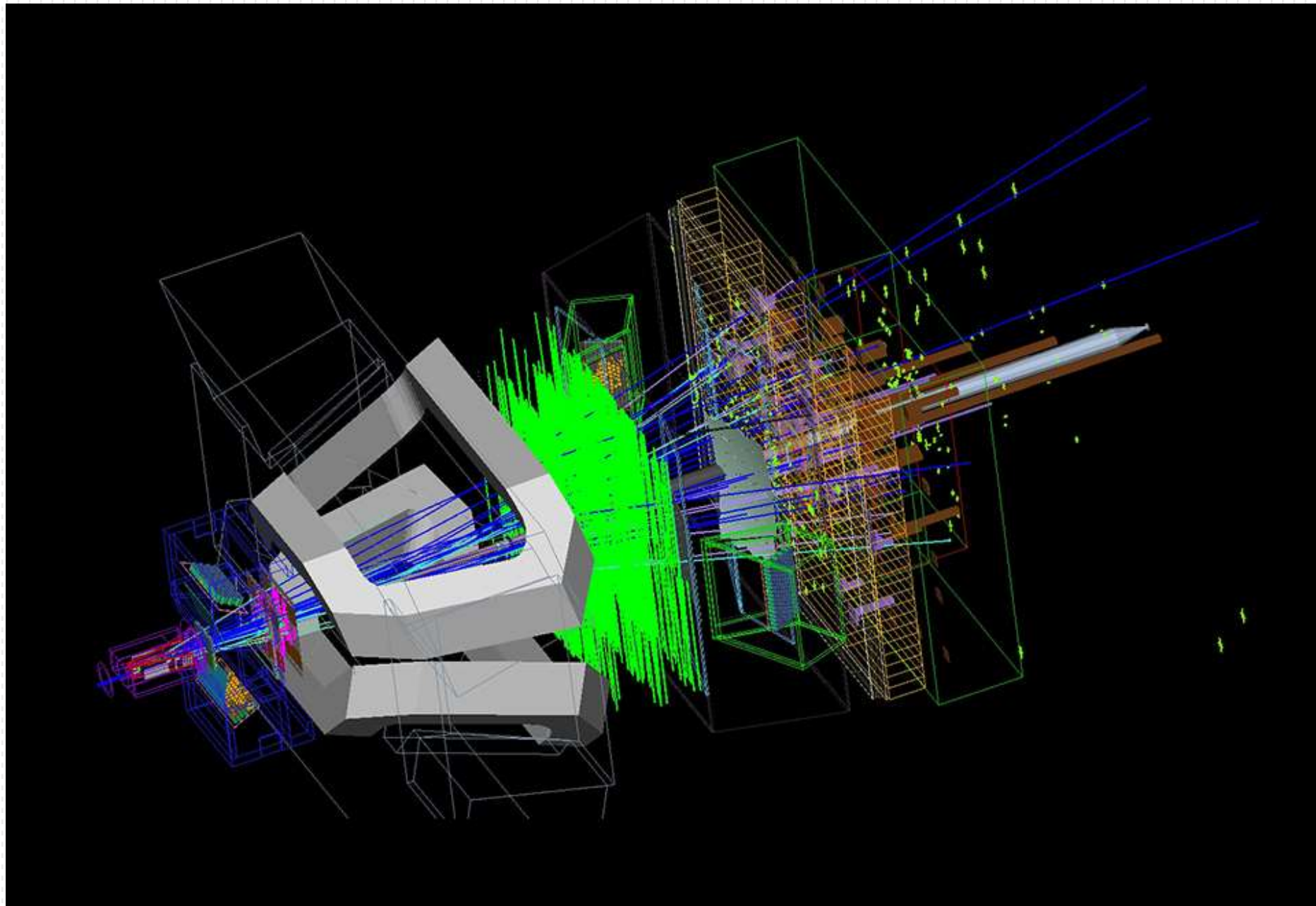


# Large Hadron Collider beauty experiment (LHCb)



### B-мезоны

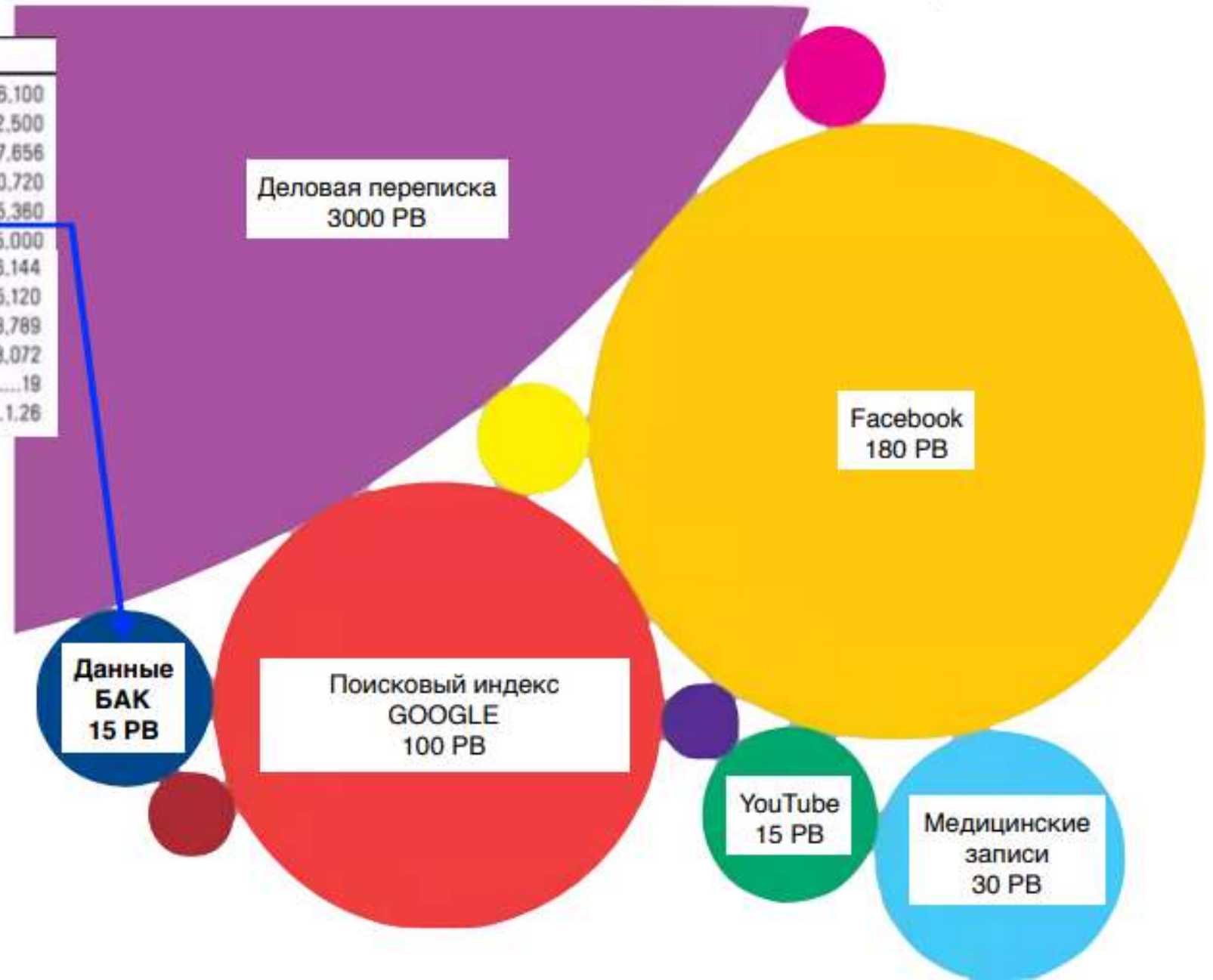
Частица ⇄	Символ ⇄	Анти-частица ⇄	Кварковый состав ⇄	Заряд ⇄	Изоспин (I) ⇄	Спин и чётность ( $J^P$ ) ⇄	Энергия покоя (МэВ/c <sup>2</sup> ) ⇄	S ⇄	C ⇄	B' ⇄	Время жизни (с) ⇄
B-мезон	$B^+$	$B^-$	$u\bar{b}$	+1	$\frac{1}{2}$	$0^-$	$5279,15 \pm 0,31$	0	0	+1	$(1,638 \pm 0,011) \cdot 10^{-12}$
Нейтральный B-мезон	$B^0$	$\bar{B}^0$	$d\bar{b}$	0	$\frac{1}{2}$	$0^-$	$5279,53 \pm 0,33$	0	0	+1	$(1,530 \pm 0,009) \cdot 10^{-12}$
Странный B-мезон	$B_s^0$	$\bar{B}_s^0$	$s\bar{b}$	0	0	$0^-$	$5366,3 \pm 0,6$	-1	0	+1	$(1,470^{+0,027}_{-0,026}) \cdot 10^{-12}$
Очарованный B-мезон	$B_c^+$	$B_c^-$	$c\bar{b}$	+1	0	$0^-$	$6276 \pm 4$	0	+1	+1	$(0,46 \pm 0,07) \cdot 10^{-12}$



Столкновение двух протонов, зарегистрированное детектором LHCb

адаптировано из <http://www.wired.com/2013/04/bigdata/>

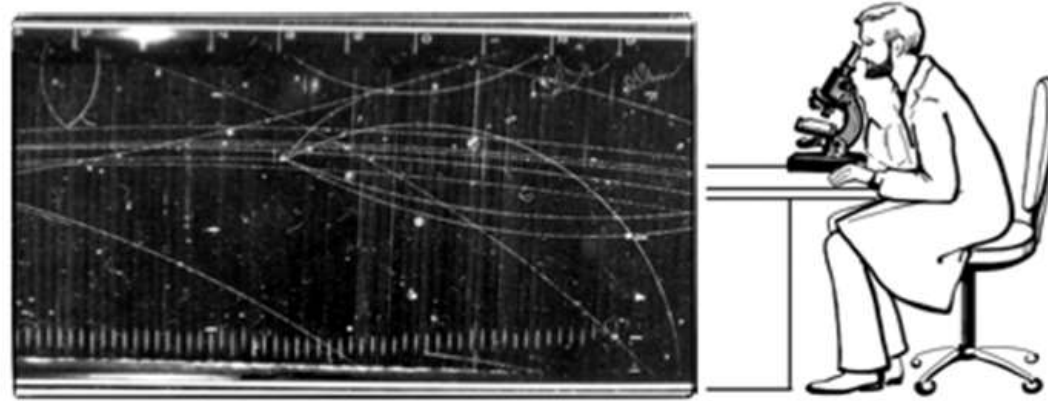
Size of data sets in terabytes	
Business email sent per year	2.986.100
Content uploaded to Facebook each year	182.500
Google's search index	97.656
Kaiser Permanente's digital health records	30.720
Large Hadron Collider's annual data output	15.360
Videos uploaded to YouTube per year	15.000
National Climactic Data Center database	6.144
Library of Congress' digital collection	5.120
US Census Bureau data	3.789
Nasdaq stock market database	3.072
Tweets sent in 2012	19
Contents of every print issue of WIRE	1.26



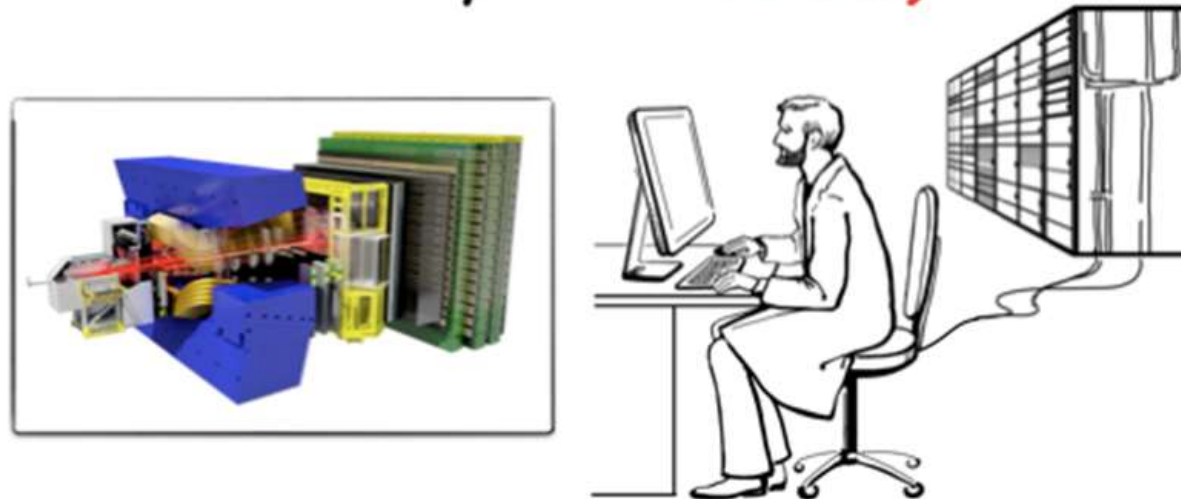
БАК создаёт  
реально большие  
данные



# What is Physics? Yesterday



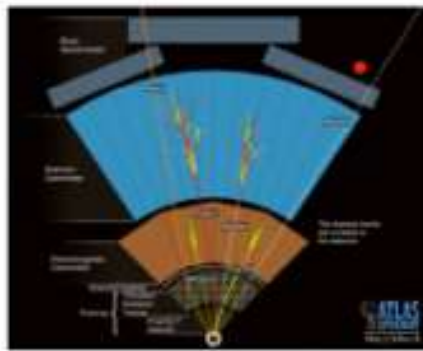
# What is Physics? Today



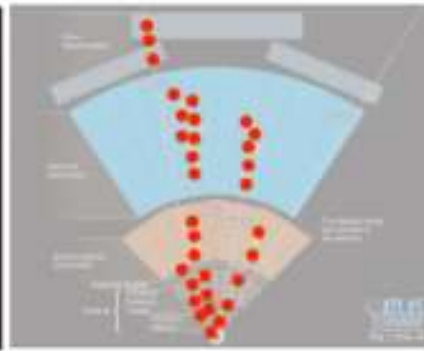
Сложные алгоритмы реконструкции событий и анализа данных

Агрессивное, многошаговое уменьшение размерности задачи

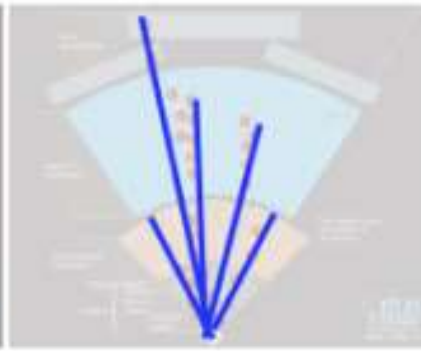
Широкий простор для применения машинного обучения и искусственного интеллекта



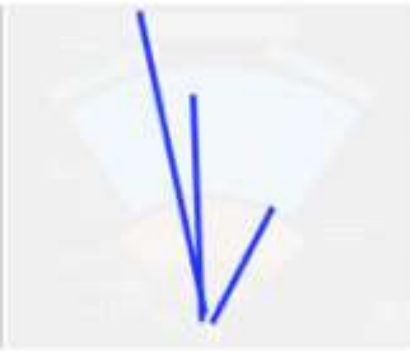
данные с  
детектора  
 $\sim 10^7$



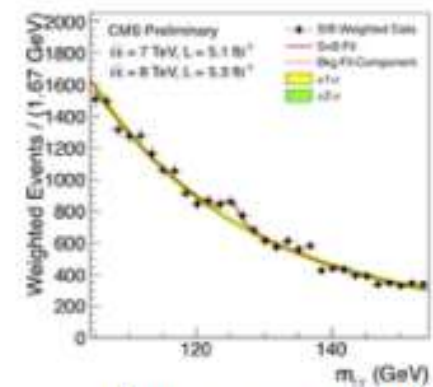
кластеризация  
 $\sim 10^4$



реконструкция  
 $\sim 50$



отбор  
 $\sim 10$



физический  
результат  
 $\sim 1$

# Yandex for CERN



Datacenter

Event search

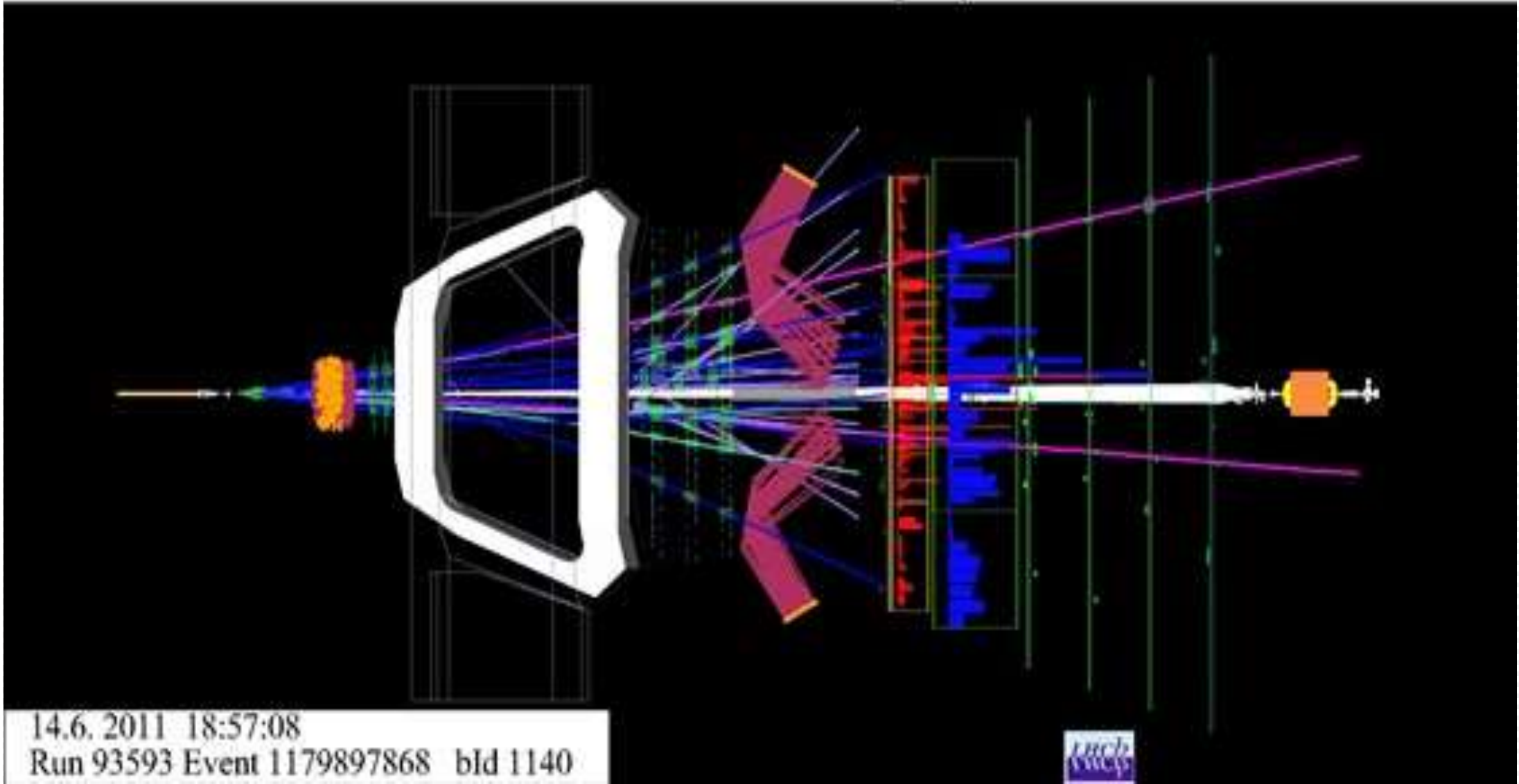
MatrixNet



Event name	Real Time Activity	Downloaded
Global Event Activity summary:		
ATLAS	2	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123
ATLAS	10	123

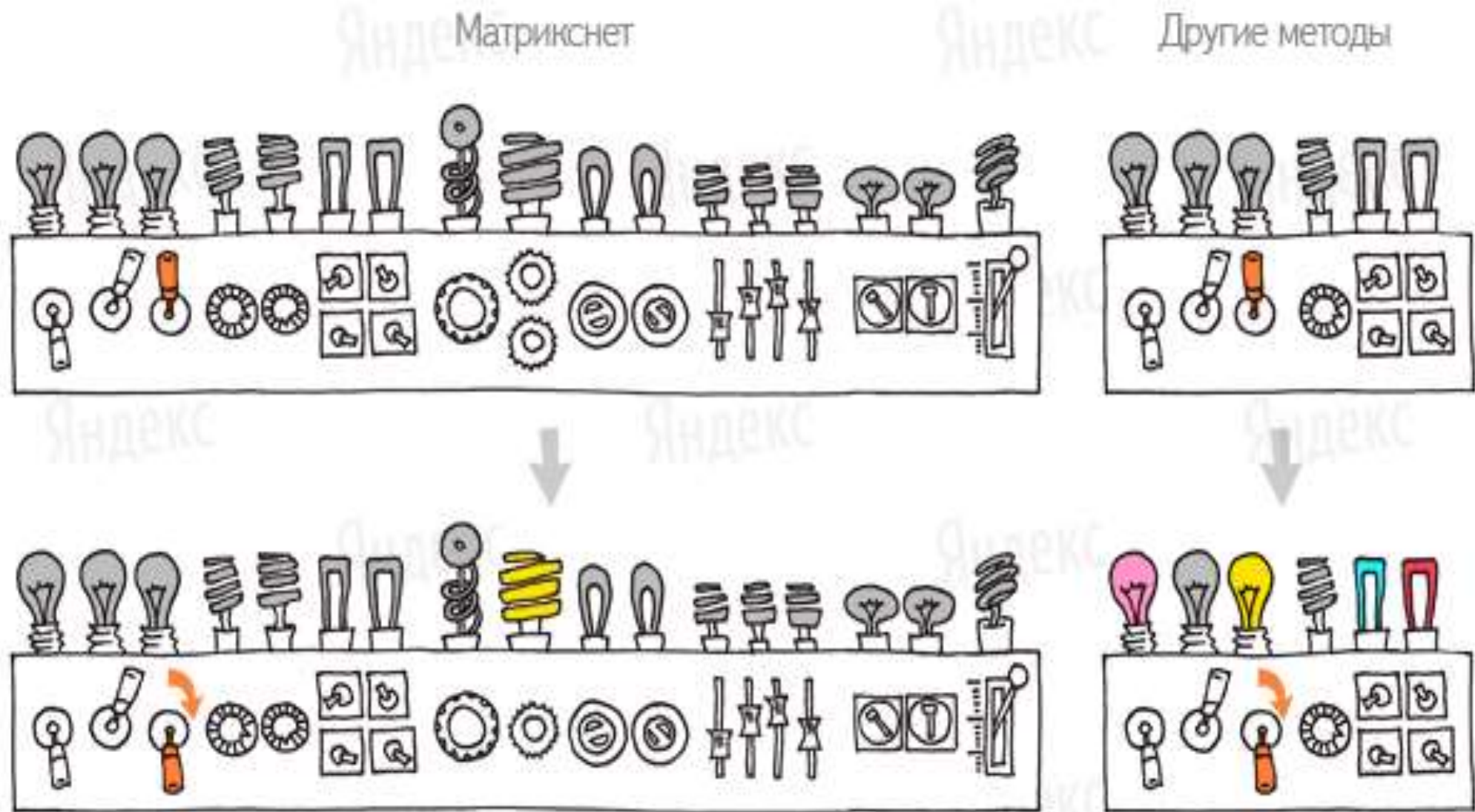


# LHCb Event Display

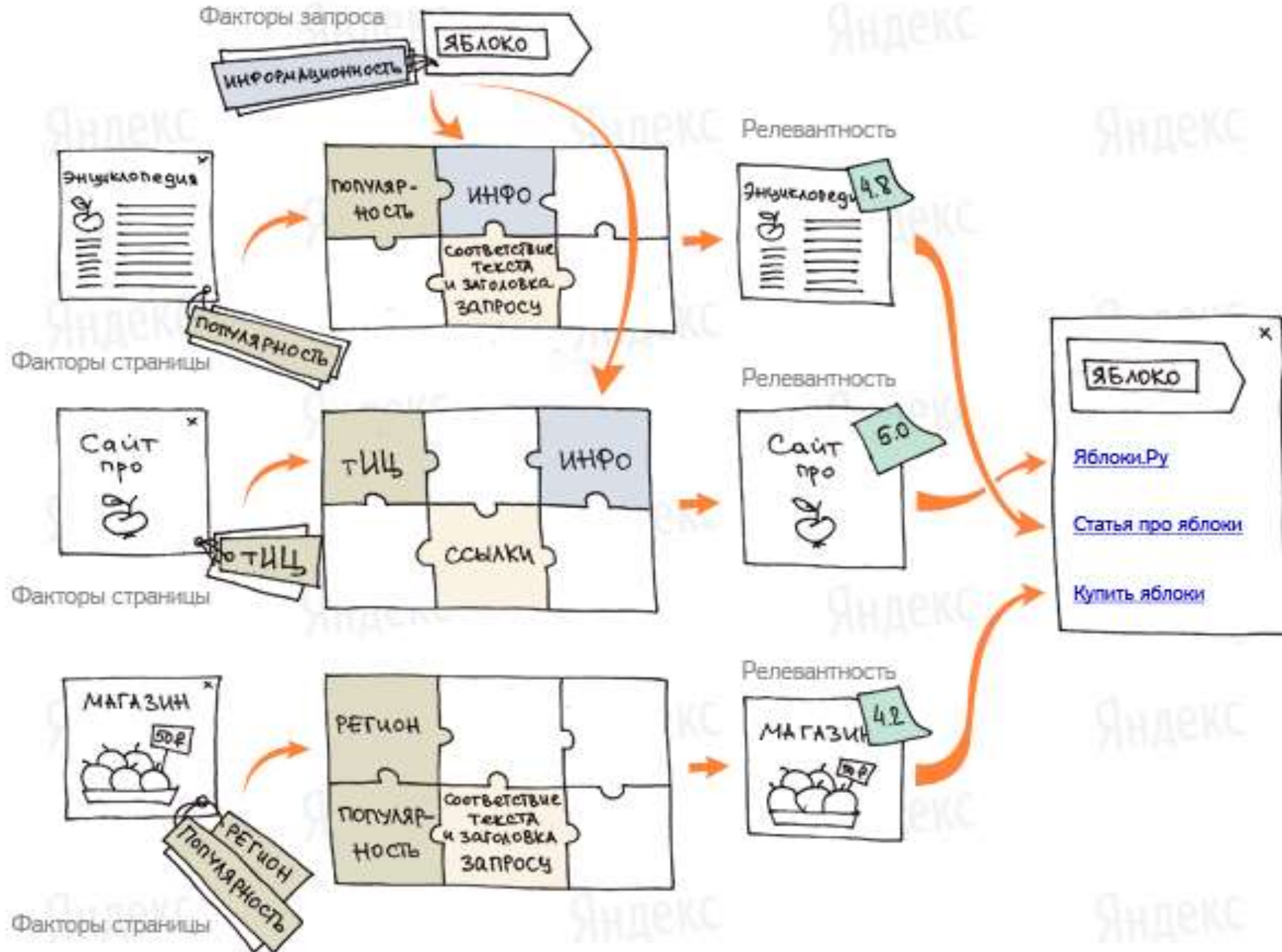


Кандидат на распад  $B_s^0 \rightarrow \mu \mu$ , наблюдаемый в эксперименте LHCb (фото ЦЕРН)

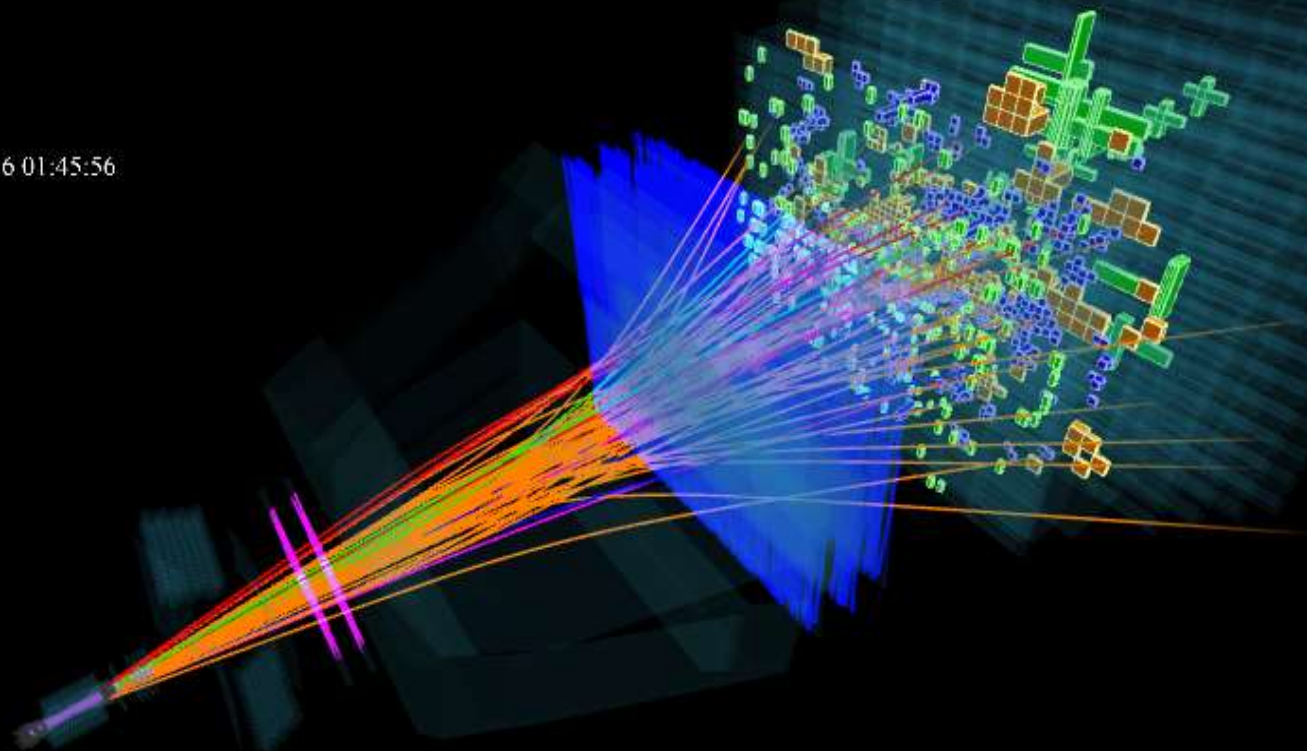
# MatrixNet



# Как устроено ранжирование



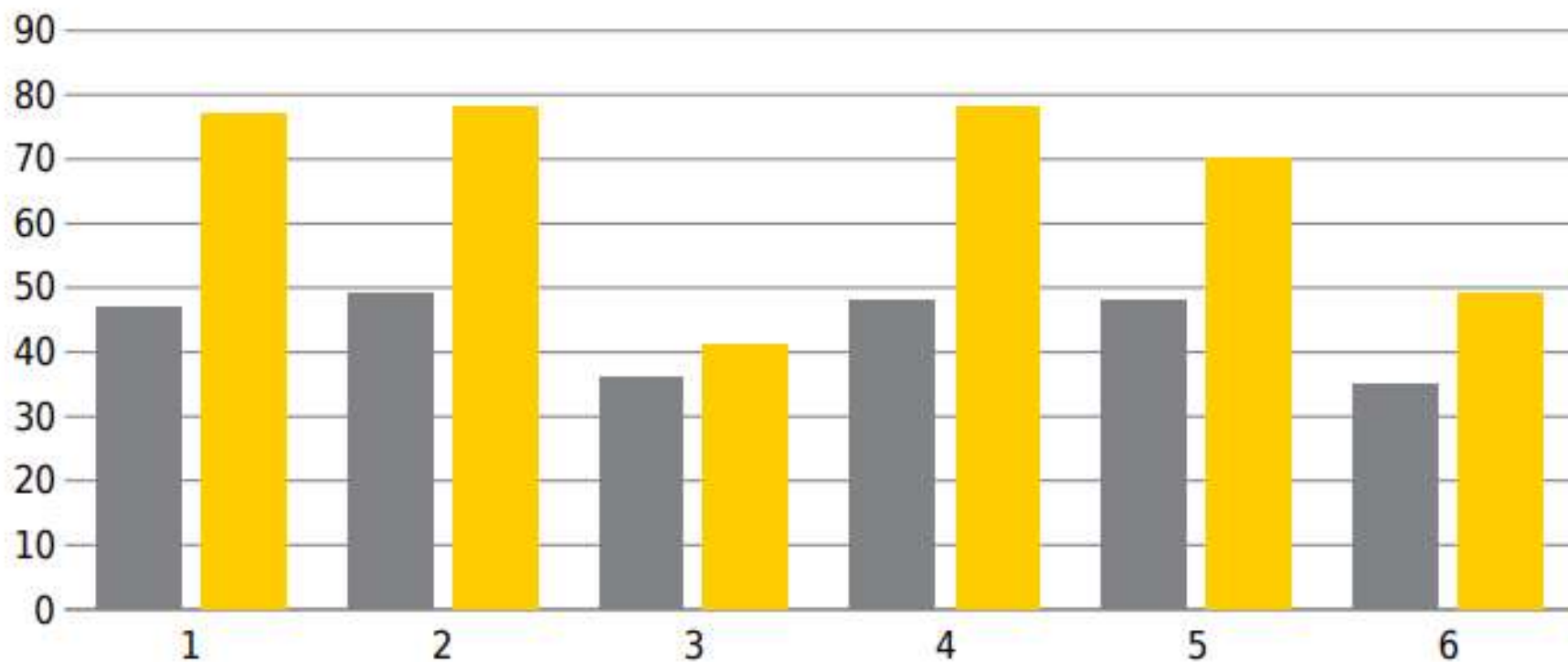
Event 74374790  
Run 173768  
Mon, 09 May 2016 01:45:56



*Типичное событие на LHCb. При столкновении протонов возник пучок частиц, которые пролетели через детектор, определивший их тип и энергию. LHCb Experiment*

## N-Body trigger Performance Comparison

(bars correspond to trigger efficiency in % for different decay modes)



■ Run-I (Before optimization) ■ MatrixNet

- ◇ Использование MatrixNet (поисковый движок Яндекса) позволило существенно улучшить эффективность триггерного отбора событий





- ◇ Данные хранятся на дисках (быстро, но дорого) и/или на лентах (дёшево, но медленно)
- ◇ С использованием машинного обучения удалось предсказать востребованность различных файлов и соответственно реорганизовать их на диске/ленте
- ◇ Результат: экономия 40% данных LHCb при ошибке предсказания ~1%



# ROOT

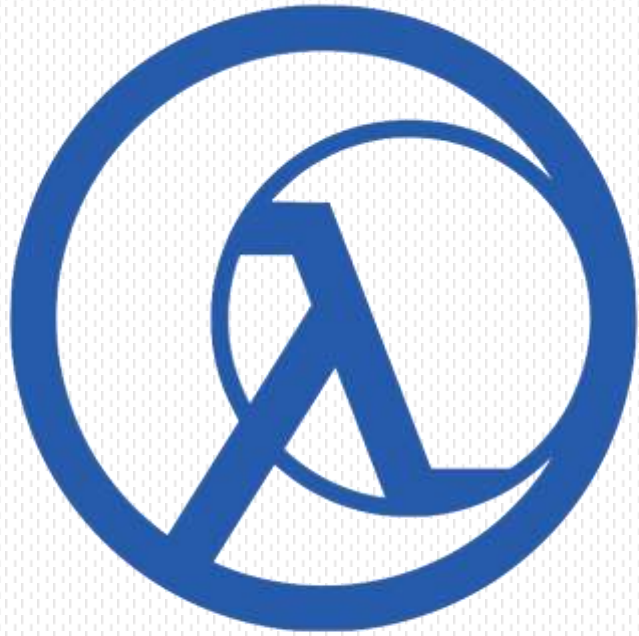
Data Analysis Framework



machine learning in Python

<http://tmva.sourceforge.net/>





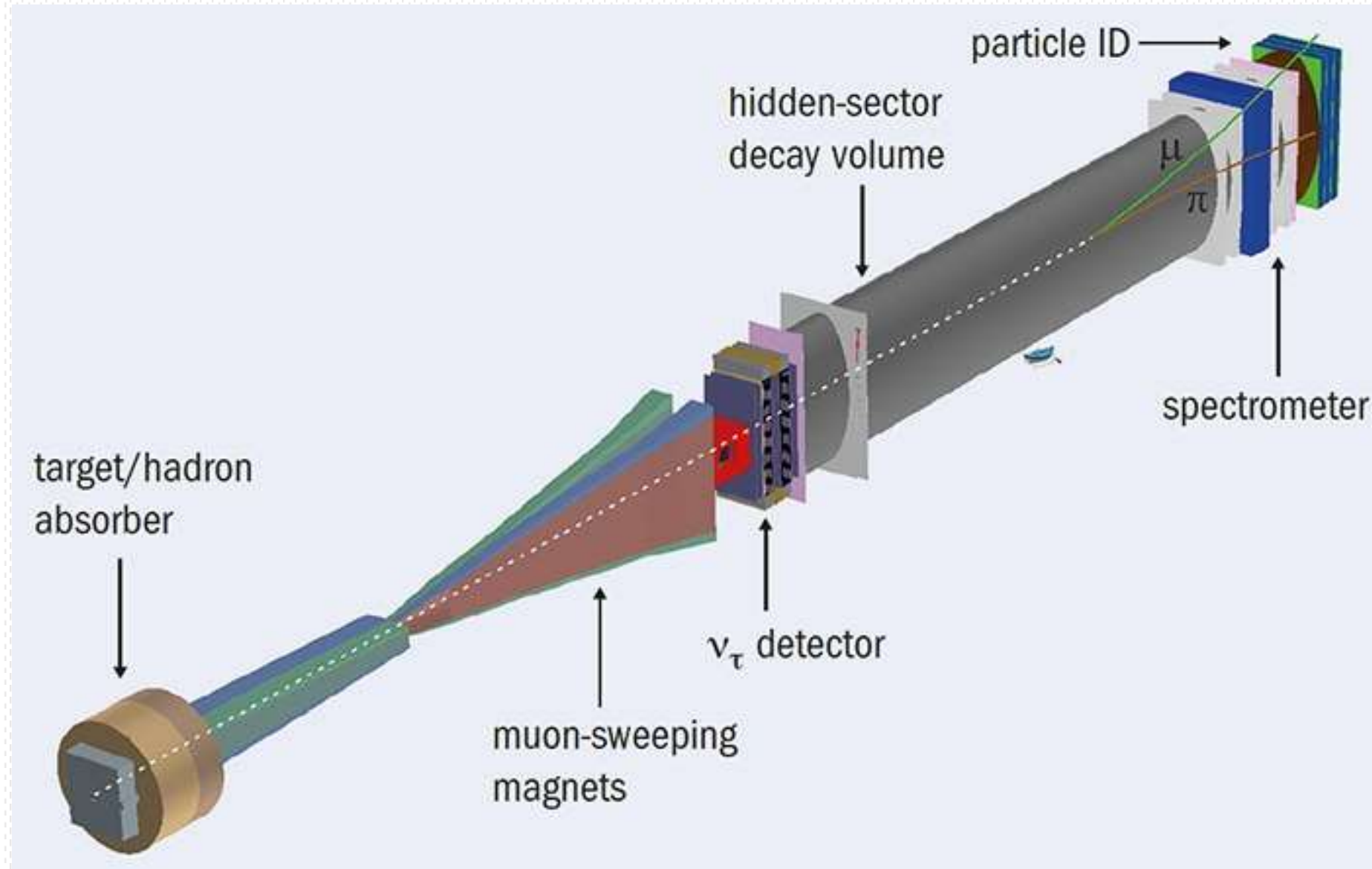
**LAMBDA • HSE**

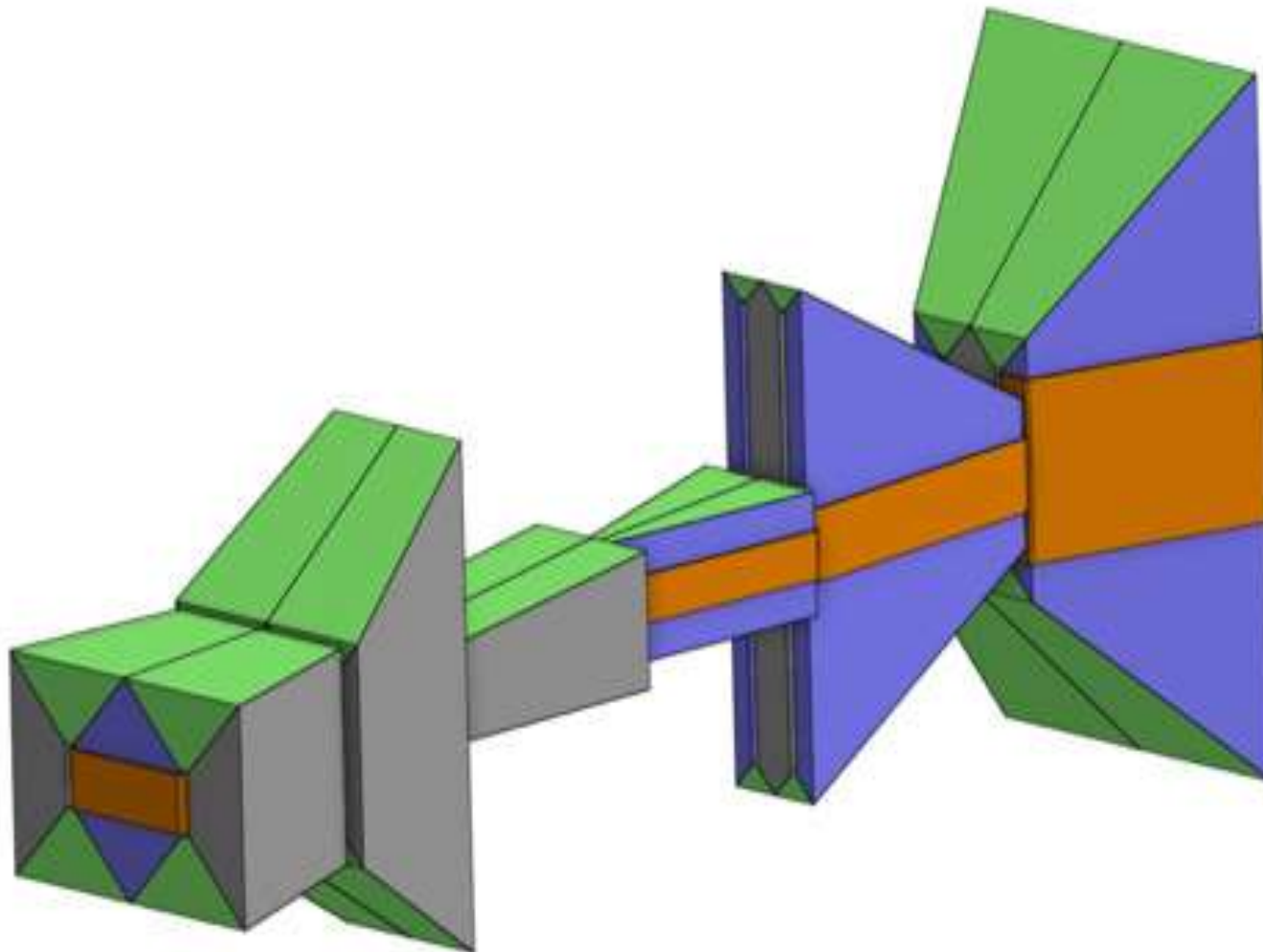
Научно-учебная лаборатория методов  
анализа больших данных



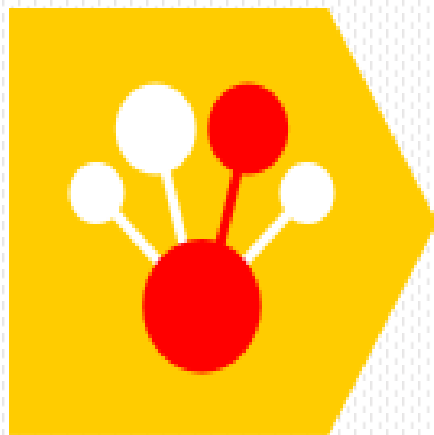
ШКОЛА АНАЛИЗА ДАННЫХ

# Общая схема эксперимента SHiP(Search for Hidden Particles)



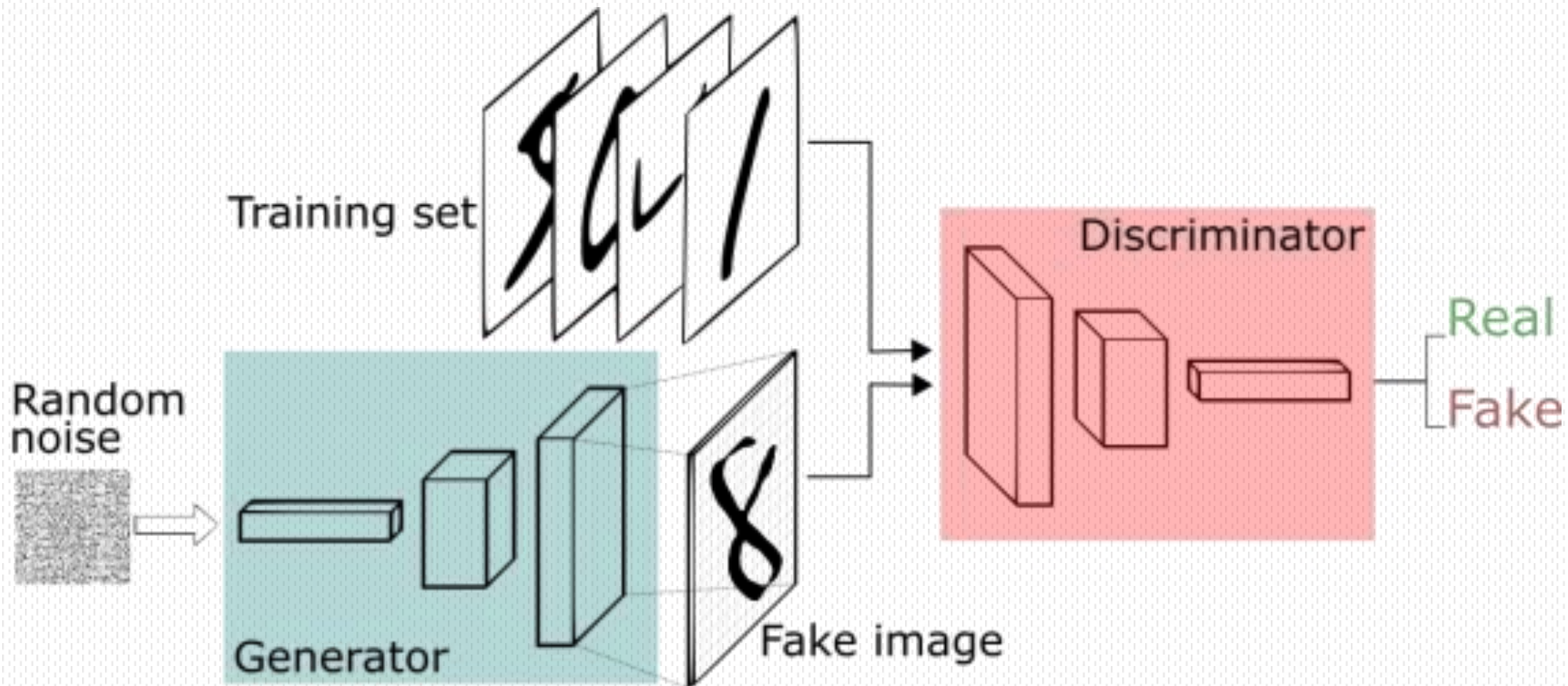


Конфигурация магнитов для эксперимента SHiP



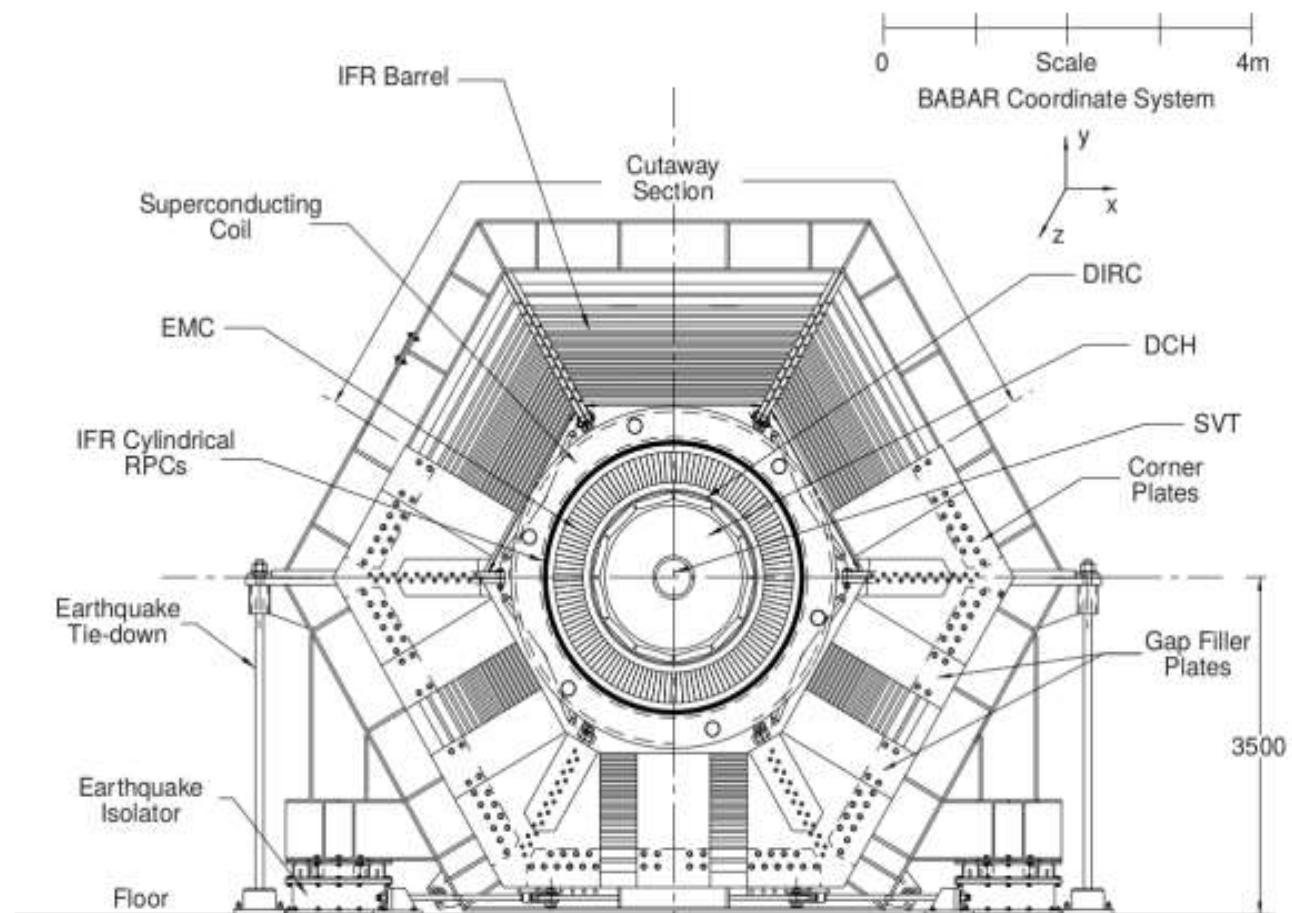
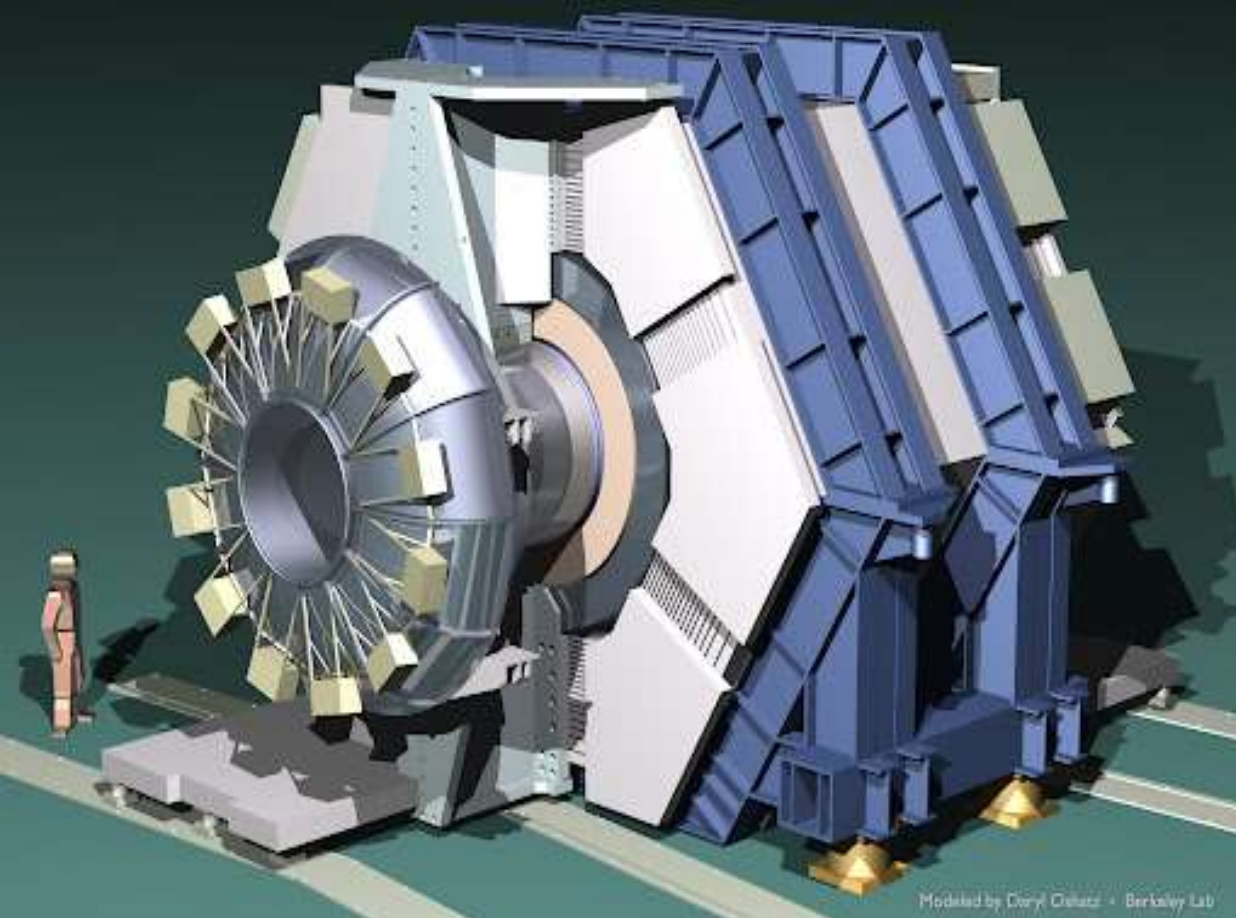
# Yandex CatBoost

**CatBoost** — открытая программная библиотека, разработанная компанией Яндекс и реализующая уникальный патентованный алгоритм построения моделей машинного обучения, использующий одну из оригинальных схем градиентного бустинга.



Генеративно-сопоставительная сеть (англ. Generative adversarial network, GAN)

# BABAR DETECTOR FOR THE PEP-II B FACTORY



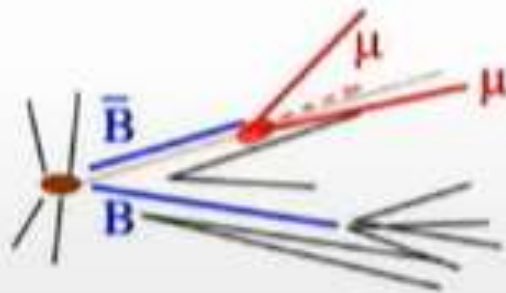
Поперечное сечение детектора *BABAR* в поперечном направлении



# Предварительный отбор (триггеры)

общее количество событий, наблюдаемых детектором за 2011-2012 гг:  $\sim 10^{15}$   
размерность каждого события  $\sim 10^7$

- частота событий на входе триггера  $\sim 10^7$  событий/секунду
- определение траекторий по хитам различных элементов детектора, размерность  $10^4$
- определение типов частиц (траекторий): мюон, каон, пион, протон, ...
- реконструкция, восстановление структуры распадов:
  - 2 трека сходятся в одну точку - *вторичная вершина*
  - *вторичная вершина* находится на расстоянии от *первичной вершины* (точки соударений протонов)
- частота событий на выходе триггера  $\sim 10^4$  событий/секунду
- сохраняются в распределенной система хранения и обработки данных



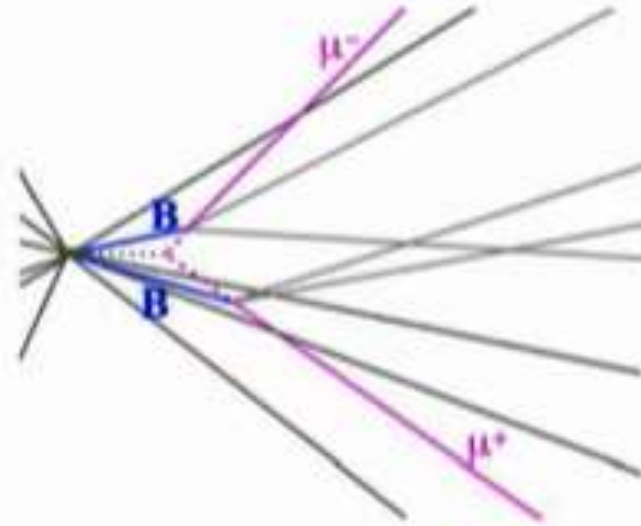
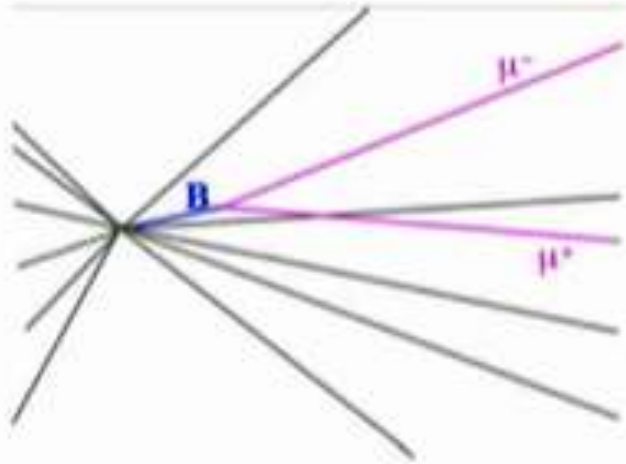
# Отбор событий-кандидатов

событие-кандидат - распад родительской частицы заданной массы на интересующие дочерние частицы. внутри одного события могут присутствовать несколько распадов

- 2 трека мюонов сходятся  $\sim$  в одну точку (вторичная вершина)
- вторичная вершина находится на расстоянии от первичной вершины
- размерность "события-кандидата" -  $10^2$
- инвариантная масса родительской частицы  $*$ , восстановленной из 2 мюонов, оказывается близкой к массам  $B_s$  или  $B_0$  ( $[4.9-6]$  ГэВ/ $c^2$ )
- каждое событие-кандидат представляем вектором признаков, размерность  $\sim 10$

$*$  вторичная вершина - точка распада частицы, породившей мюоны. поэтому эту частицу называют *родительской частицей*

# Сигнал, фон



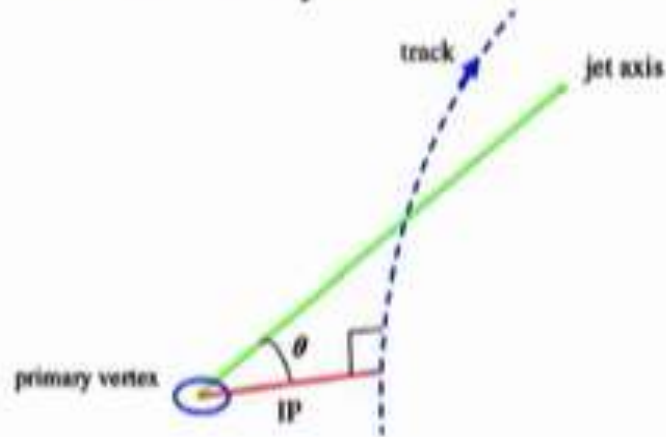
## Сигнал:

- два мюона
- инвариантная масса в районе массы  $B_s$ ,  $B^0$
- хорошо реконструированная вторичная вершина на заметном ( $\sim 2$ мм) расстоянии от первичной вершины, импульс которой сонаправлен с направлением полета

## Фон:

- два распада  $B$  на  $\mu$  и  $X$
- распад  $B$  на  $\mu$  и ошибочно идентифицированный мюон
  - $B_s \rightarrow K^- \mu^+ \nu$
- Одиночные распады  $B$ :
  - $B_s \rightarrow K^+ K^-$
- ...

# Признаки событий, В-мезон

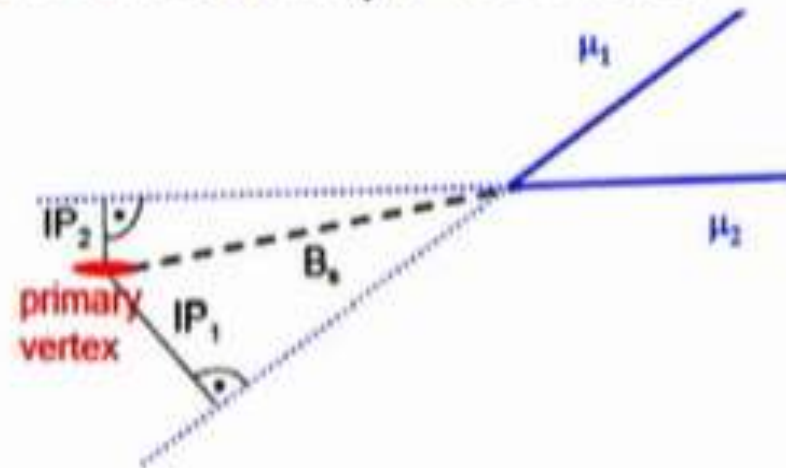


точка распада В - вторичная вершина распада

1. время жизни В
2. IP - расстояние от траектории В до первичной вершины
3.  $p_T$  - поперечный импульс в системе координат детектора
4. изоляция В
5. угол между импульсом В и  $P_{thrust}$
6. угол между импульсом  $\mu^+$  в системе покоя В и  $P_{thrust}$  в системе покоя В

$P_{thrust}$  - сумма импульсов всех треков, начинающихся из точки распада другого b-адрона

# Признаки событий, мюоны

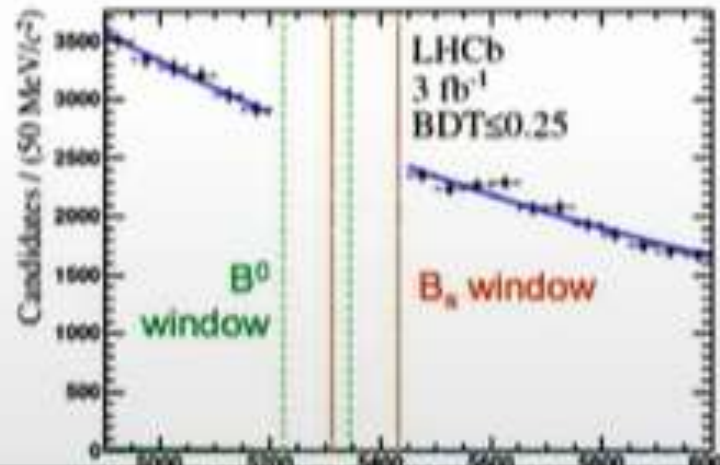


1. минимальная значимость прицельного параметра
2. минимальное расстояние между 2 мюонами
3. изоляция трека
4. угол поляризации (the cosine of the angle between the muon momentum in the dimuon rest frame and the vector perpendicular to both the  $B$  candidate momentum and the beam axis)
5. модуль разницы между углами  $\phi$  мюонов в сферической системе координат
6. модуль разницы между псевдобыстротой мюонов

# Сигнальные регион, боковой регион, слепой анализ

- ожидаем сигнал в районе инвариантной массы родительской частицы
  - комбинаторный фон
  - сигнал (*hopefully*)
  - "пиковый" фон с двойной ошибкой идентификации
- вырезаем сигнальный регион и прячем до лучших времен
- боковые регионы - содержат только фон (различный)

**ВАЖНО:** ожидаем, что распределение комбинаторного фона по массе должно принадлежать одному классу ( $e^{-m}$ ) в сигнальном и боковом регионах

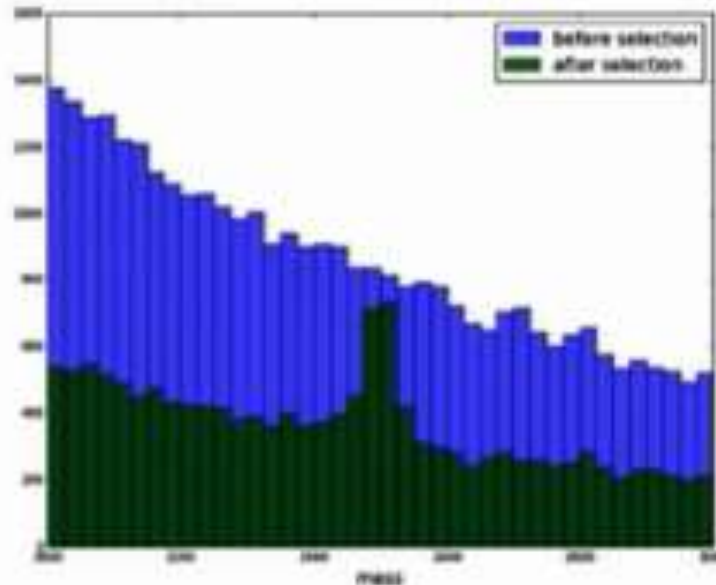


# Использование машинного обучения

- *не можем разметить сигнальные события в реальных данных!*
- сгенерировать распады для сигнальных каналов
  - конфигурация распада элементарных частиц
  - отклик детектора на распады
- сгенерировать распады для фона
- провести сгенерированные данные через те же этапы отбора событий как и реальные события
  - триггеры
  - идентификация частиц
  - реконструкция
- привести к признакам описания событий-кандидатов
- датасет из 12 колонок (признаков) с разметкой:
  - 1 - сигнал
  - 0 - фон

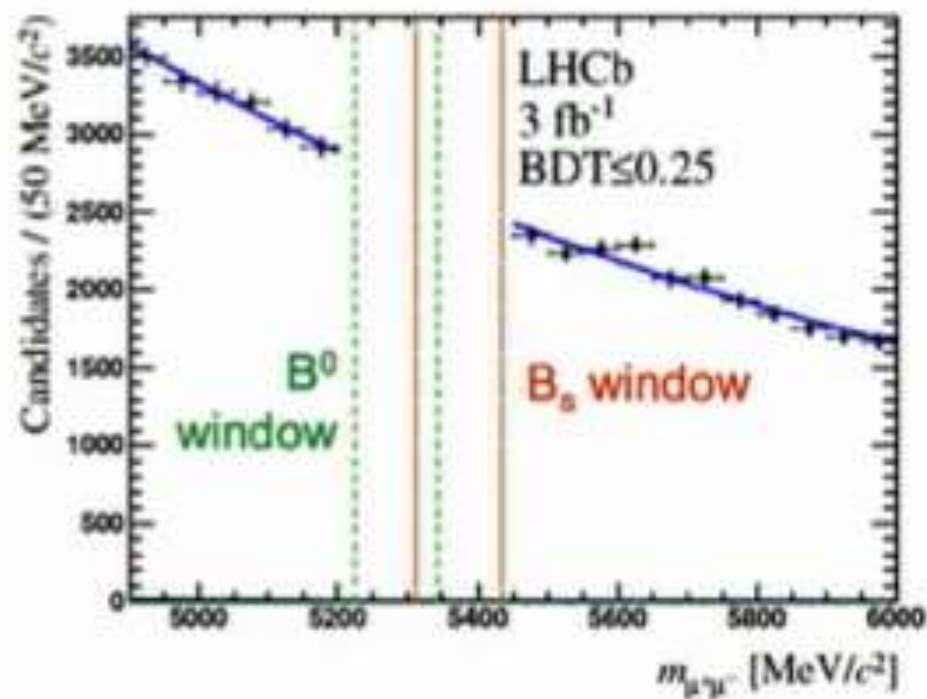
# Особенности обучения

- результат работы классификатора - оценка *количества* сигнальных событий, а не вероятности принадлежности отдельного события выбранному каналу
- другими словами, с помощью классификатора  $g$  мы хотим найти область  $G : g(x) > p$ , в которых находится как можно больше сигнальных (true positive) событий и как можно меньше фоновых (false positive)
- **распределение фоновых событий, отобранных классификатором в сигнальном и боковых регионах, должно принадлежать одному классу, чтобы не создавать "искусственных пиков":**



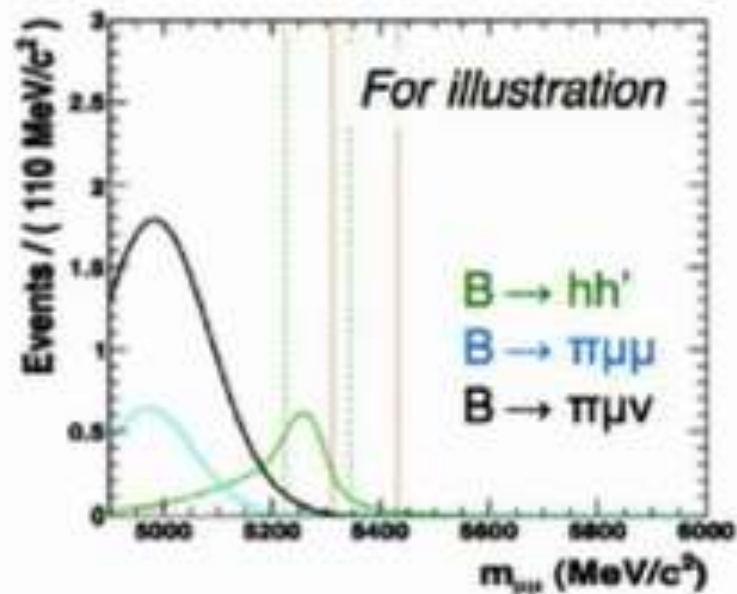


# Комбинаторный фон



- основной источник фона распадов вида  $b\bar{b} \rightarrow \mu\mu X$
- оцениваем ожидаемое количество фоновых событий в сигнальном регионе, аппроксимируя экспоненциальное распределение комбинаторного фона, на сигнальный регион для каждого интервала на значения порогов классификатора
- при высоких значениях на порог классификатора практически исчезает

# Пиковые фоны

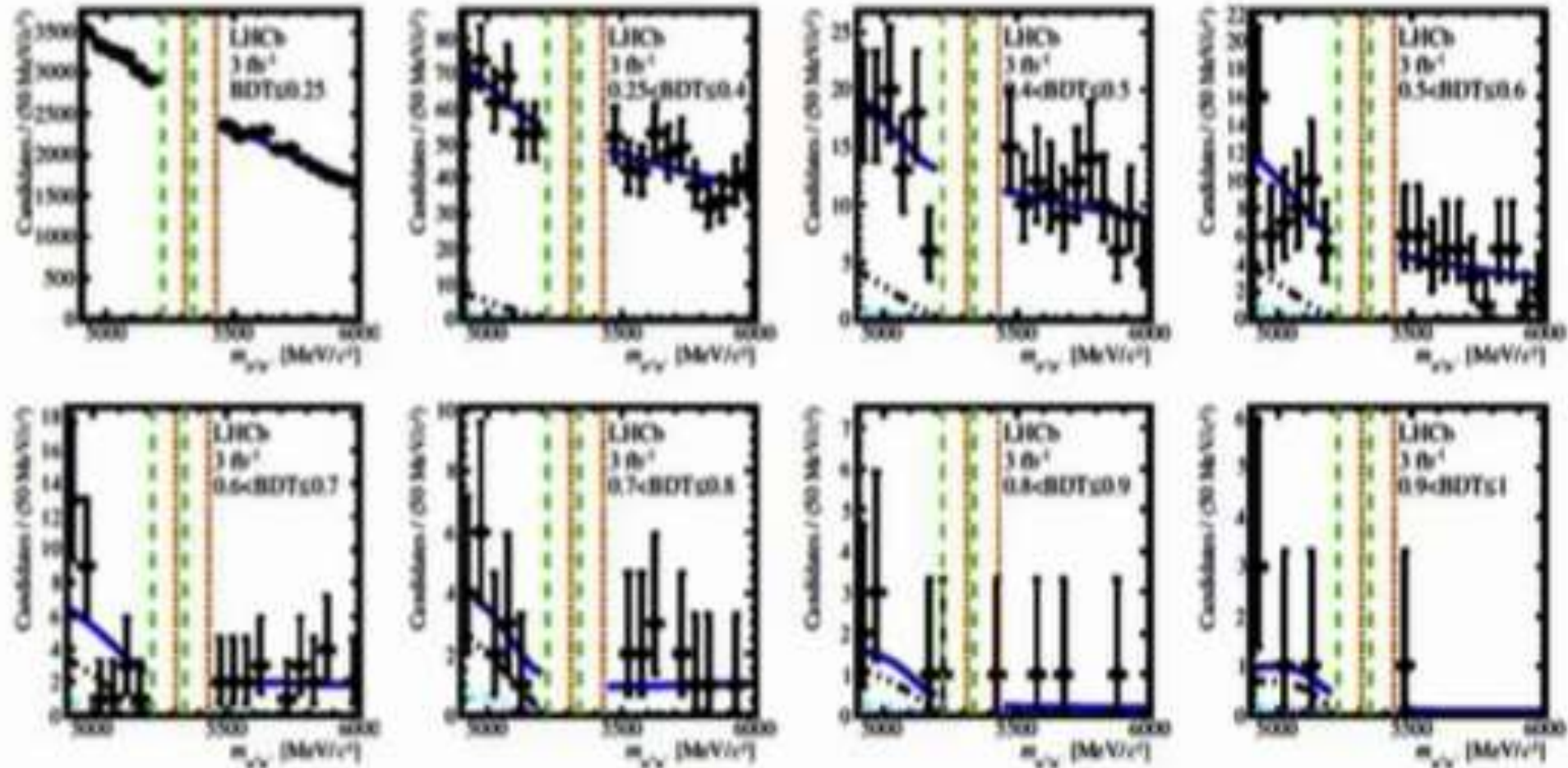


- в сигнальном регионе играет роль только  $B \rightarrow h^+ h'^-$  (двойная мисидентификация мюонов)
- в боковых регионах, распады с одиночной ошибкой идентификации или 2 мюонами исходящими из одной вершины:
- $B^0 \rightarrow \pi^- \mu^+ \nu$ ,  $B^{0/+} \rightarrow \pi^{0/+} \mu\mu$
- $B_s \rightarrow K^- \mu^+ \nu$ ,  $B_c \rightarrow \pi^{0/+} J/\psi(\mu\mu)\mu\nu$
- $\Lambda_b \rightarrow p\mu\nu$
- прочие каналы вносят несущественный вклад

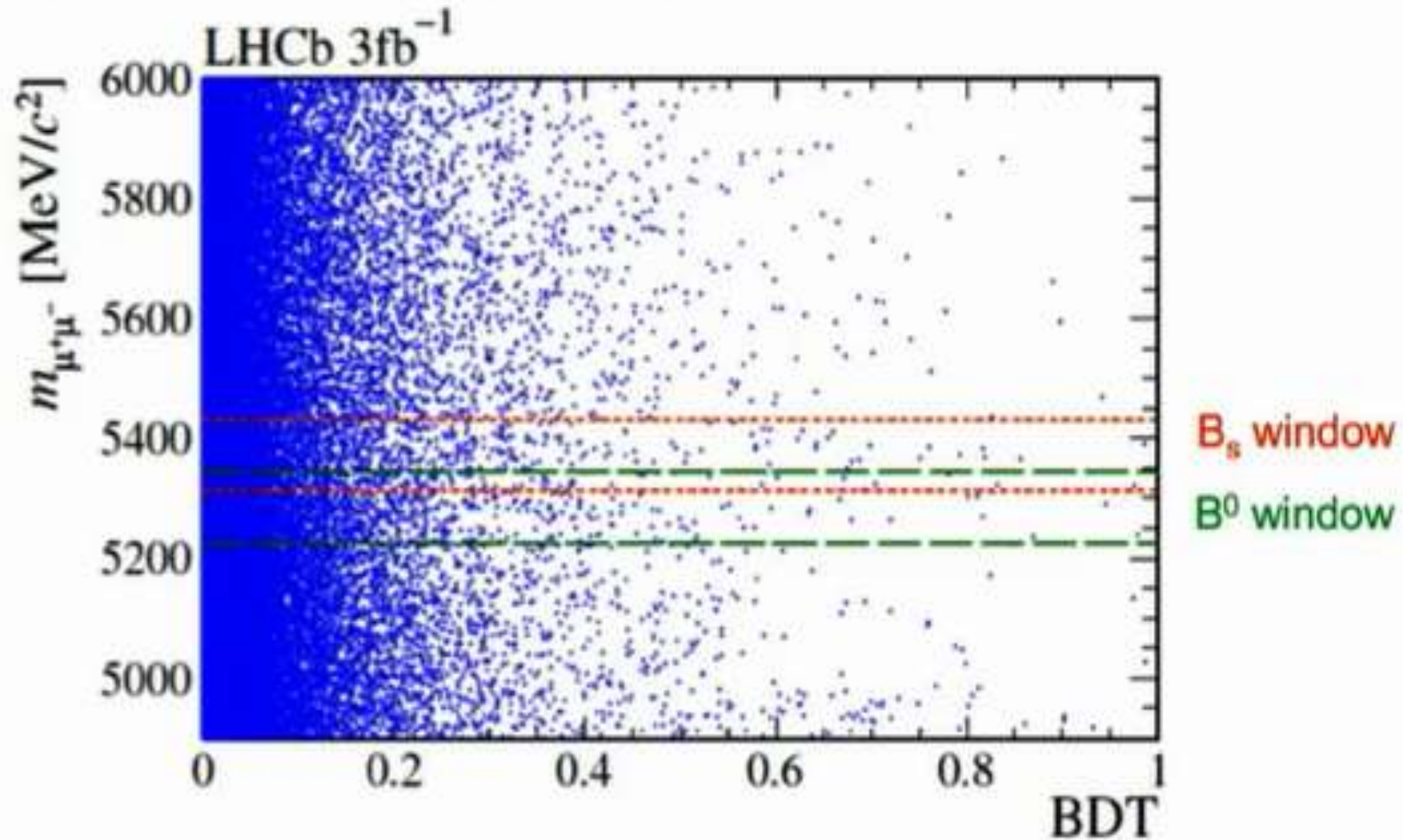
# Аппроксимация фоновых каналов

$B^0 \rightarrow \pi \mu^* \nu$ ,  
 $B_s \rightarrow K \mu^* \nu$ ,  
 $B^{0*} \rightarrow \pi^{0*} \mu \mu$

$B_{d/s} \rightarrow h^* h'$   
 total



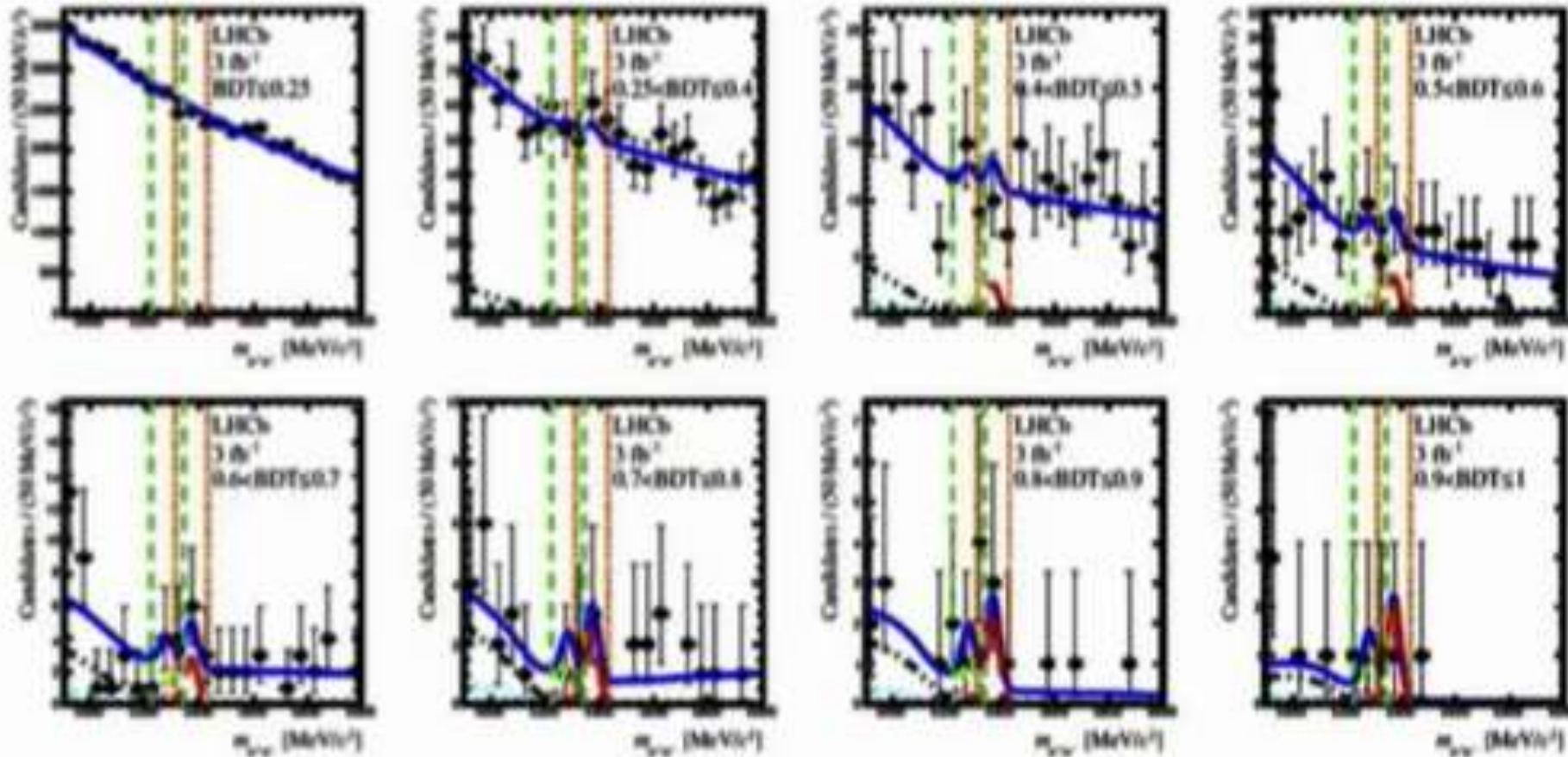
# Результаты применения классификатора к реальным данным



# Аппроксимация сигнального региона

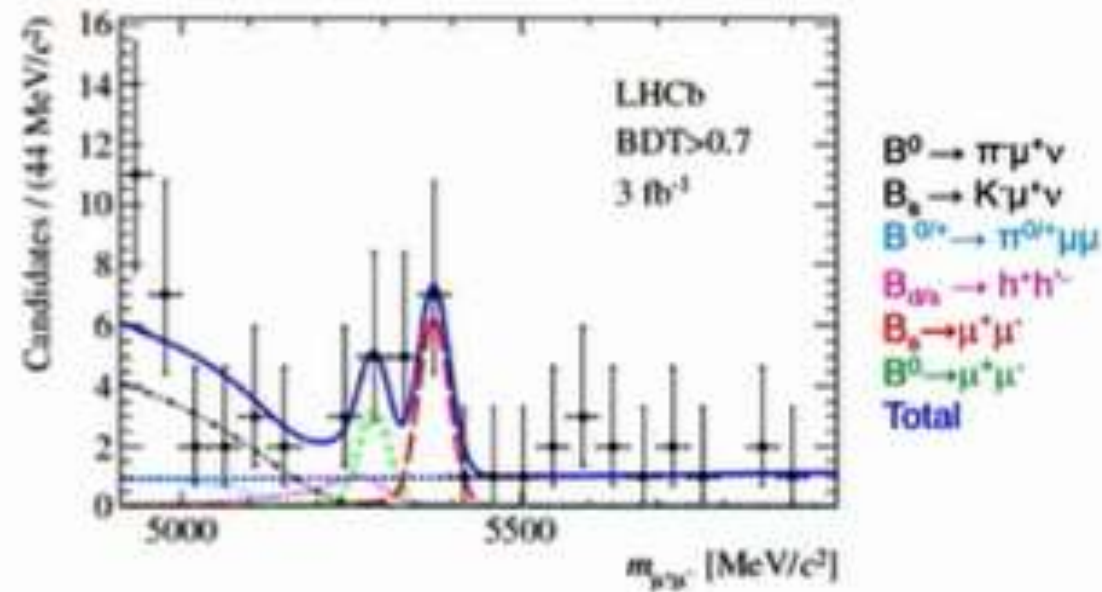
$B^0 \rightarrow \pi \mu^+ \nu$   
 $B_s \rightarrow K \mu^+ \nu$   
 $B^{0*} \rightarrow \pi^{0*} \mu \mu$

$B_{\psi_s} \rightarrow h^+ h^-$   
 $B_s \rightarrow \mu^+ \mu^-$   
 $B^0 \rightarrow \mu^+ \mu^-$   
**Total**



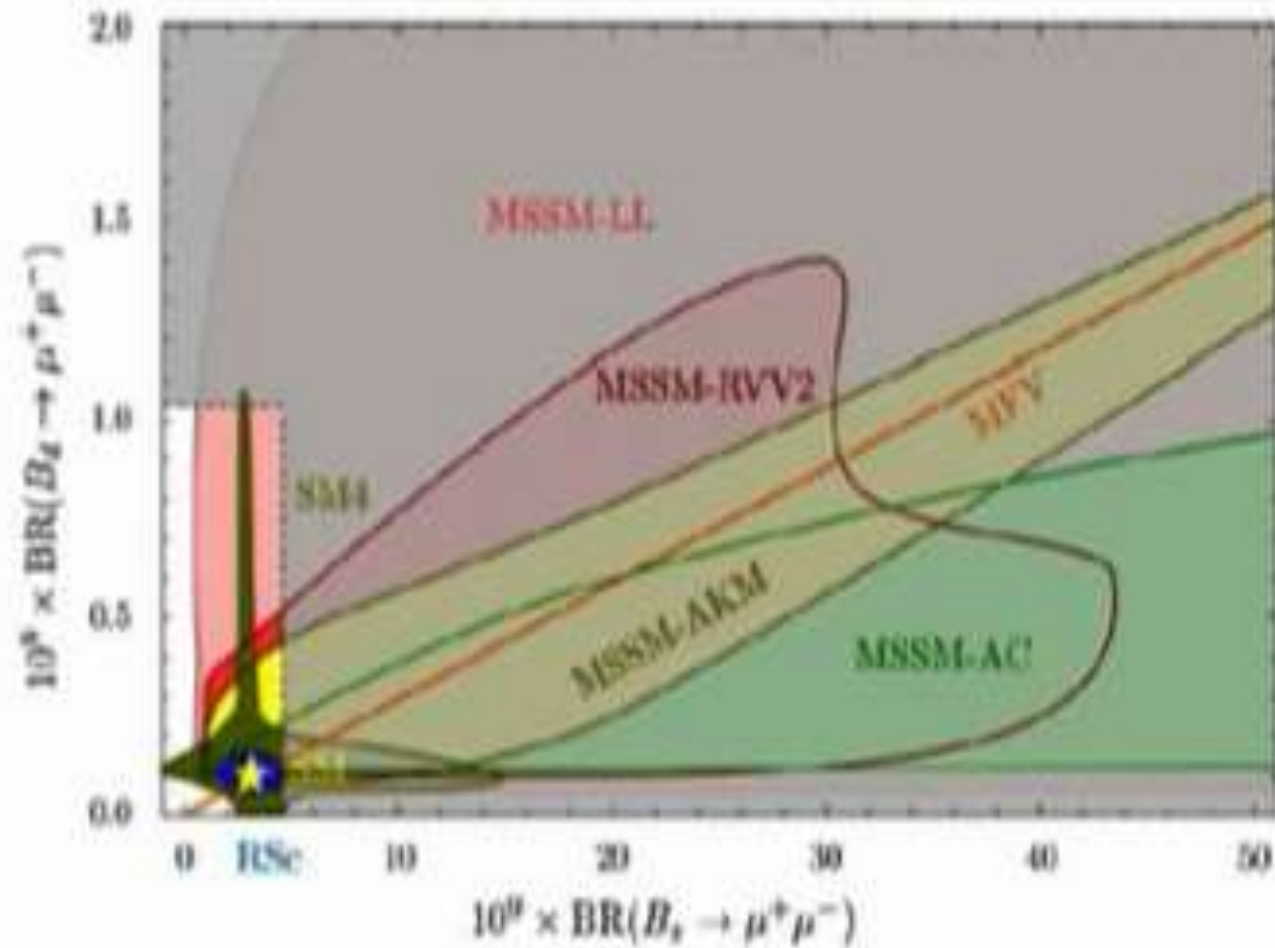
# Подсчет событий

подсчет событий для области  $G : g(x) > 0.7$



- $BF(B_s \rightarrow \mu\mu) = (2.9_{-1.1}^{+1.4})10^{-9} (4\sigma)$
- $BF(B^0 \rightarrow \mu\mu) = (3.7_{-2.1}^{+2.4})10^{-10} (2\sigma)$
- Предсказания стандартной модели:
  - $BF(B_s \rightarrow \mu\mu) = (3.25 \pm 0.17)10^{-9}$
  - $BF(B^0 \rightarrow \mu\mu) = (1.07 \pm 0.10)10^{-10}$

# Влияние на новую физику (2012)



**Спасибо за внимание!**



