# <span style="color:red">Сбор</span> открытых данных

Москва, 2016

# Проект Import.io

# Google Spreadsheets

# OpenRefine

# Outwit Hub

# Scrapy

@iradche
http://www.slideshare.net/iradche
https://www.facebook.com/iradche
iradche@gmail.com

# Ирина Радченко

Кандидат технических наук,
Доцент университета ИТМО
Сооснователь проекта DataDrivenJournalism.Ru
Основатель ODI Moscow (позднее - ODI St.Petersburg)
Координатор OKI-RU