

The Yandex Leak: How a Russian Search Giant Uses Consumer Data

Kaileigh McCrea, Privacy Engineer, Confiant

About Me

Kaileigh McCrea

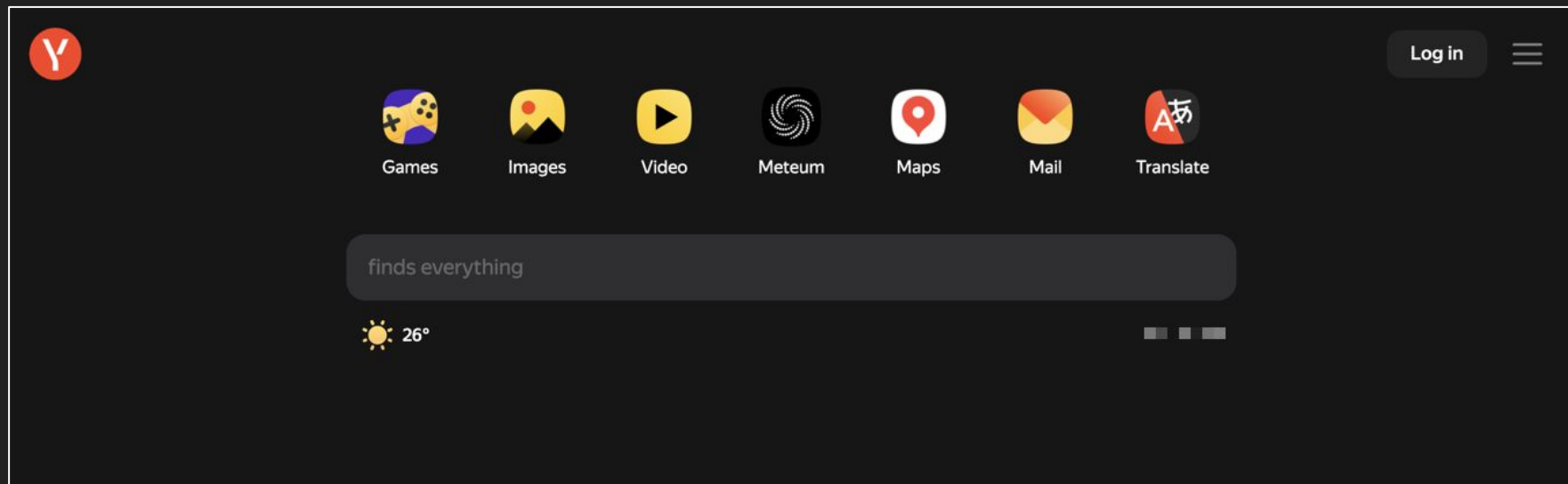
- Privacy Engineer at Confiant (3 yrs)
- Software Engineer (6 years)
- Cybersecurity Nerd
- Recovering Political Science major
- Twitter: @kaileighrose

Roadmap

- Background on Yandex Leak
- Dive into code:
 - What data Yandex is collecting
 - What Yandex is doing with that data
 - Who Yandex is sharing that data with
- Conclusions and wrap up
- Q&A

Yandex 101

Yandex 101



Yandex 101

Key Businesses



Yandex 101



AppMetrica: “In-depth analytics for product and growth teams”

Crypta: “helps to identify important user characteristics for advertisers”



Audiences: allows you to pull data from several sources to generate your own targeted segments

Yandex 101

Yandex LLC

Head office in Russia: Moscow

Head office

16, Leo Tolstoy St., Moscow, Russia
119021
tel.: +7 495 739-70-00
fax: +7 495 739-70-70

Advertising clients

tel.: +7 495 739-37-77
fax: +7 495 739-23-32
adv@yandex-team.ru

Investor Relations

tel.: +7 495 974-35-38
askIR@yandex-team.ru

Public relations

pr@yandex-team.ru

Corporate Secretary

secretary@yandex-team.ru

Sustainability

sustainability@yandex-team.com

Official Telegram channel for individual investors https://t.me/yndx_forinvestors (in Russian only)

Yandex N.V.

Registered office in Amsterdam

Schiphol Boulevard 165, 1118 BG Schiphol, The
Netherlands
tel.: +31 0 20 206 6970

Yandex: A Drama

Russian internet giant grants veto powers to Kremlin-linked body

Yandex agrees to corporate restructuring in move likely to increase government oversight

Andrew Roth in Moscow

Mon 18 Nov 2019 06.30 EST



Arkady Volozh, the chief executive of Yandex, said the company would maintain control over its daily operations. Photograph: Mikhail Metzel/Tass

Warnings raised over Russian tech giant Yandex's UK operation

MPs want restrictions placed on the company, known as Russia's Google, which also runs the Yango Deli grocery service

● [Russia-Ukraine war: live news](#)

Shanti Das

Sat 5 Mar 2022 15.02 EST



📷 A Yango Deli driver on an electric moped delivers to homes in London. The service is expanding across the city. Photograph: John Sibley/Reuters

Data-harvesting code in mobile apps sends user data to “Russia’s Google”

Data from apps on Apple- and Google-powered mobile devices is sent to Russian servers.

PATRICK MCGEE, FINANCIAL TIMES - 3/29/2022, 7:18 AM



Yandex

User profiles

AppMetrica: Your app’s CRM

Build complete audience knowledge with segmentation based on profile data or dive into individual users with profile cards.

Today

- Launch app
- Start onboarding
- Go to catalog
- View item
- Add to cart

Russia's war hits Yandex, the 'Google of Russia'

Sources say the company is seeking a media exit as top exec hit with sanctions over propaganda charge

Natasha Lomas, Ingrid Lunden / 12:20 PM PDT • March 16, 2022

 Comment



[Premium](#) [HOME](#) > [TECH](#)

'I bought a plane ticket and left 12 hours later': Engineers at Yandex, Russia's Google rival, are fleeing abroad and leaving spouses and salaries behind

Rosie Bradbury Apr 12, 2022, 3:35 AM PDT

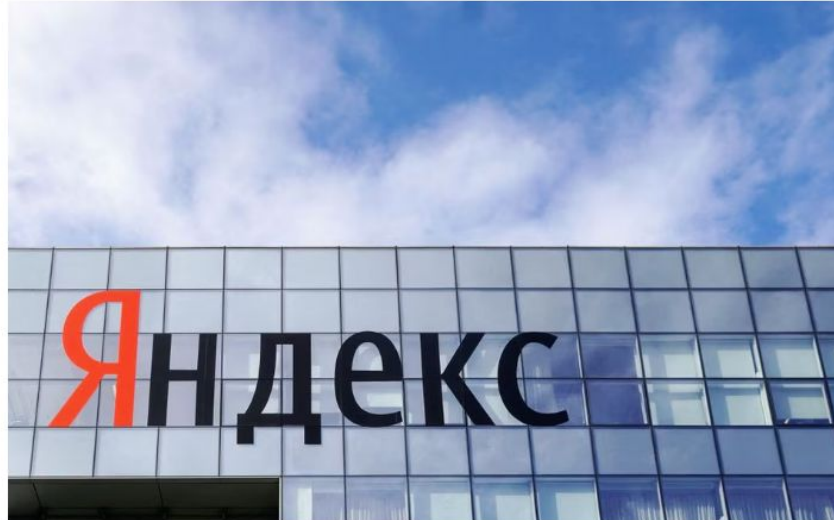


Technology

Yandex CEO resigns after being targeted by EU sanctions

Reuters

June 3, 2022 7:35 AM PDT · Updated a year ago



The logo of Russian internet group Yandex is pictured at the company's headquarter in Moscow, Russia October 4, 2018. REUTERS/Shamil Zhumatov

June 3 (Reuters) - Russian internet giant Yandex ([YNDX.O](#)) said on Friday that Arkady Volozh had stepped down as CEO and left the board of directors after the European Union included him on its latest list of sanctions against Russian entities and individuals.



Join TechCrunch+

Login

Yandex's sale of media assets to VK includes yandex.ru homepage

Natasha Lomas @riptari / 12:05 AM PDT • August 23, 2022

 Comment

Search Q

TechCrunch+

Startups

Venture

Security

AI

Crypto

Apps

Events

More



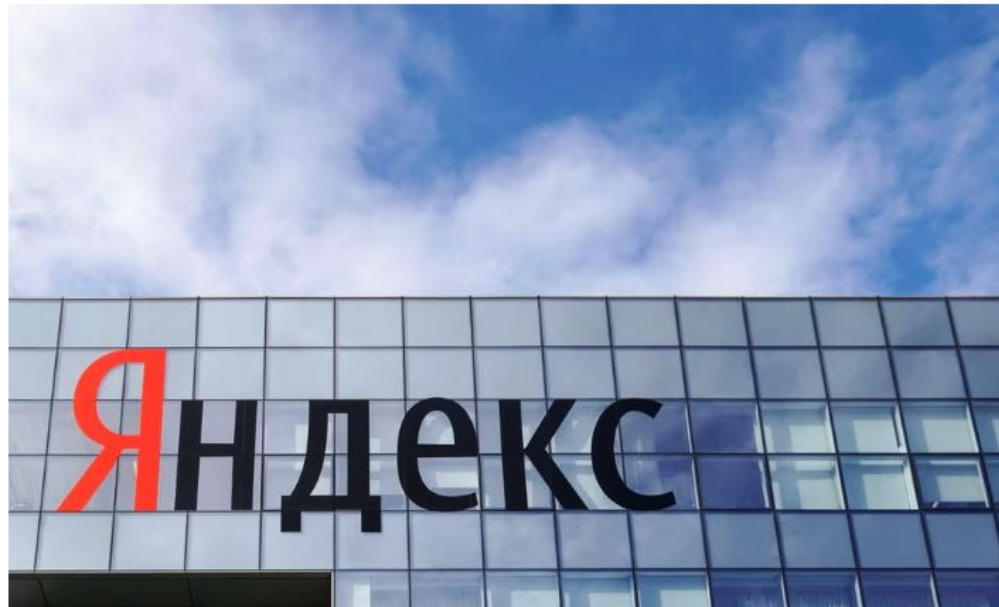


Deals

Yandex parent to review ownership of Russian tech giant, seek divestment

By Alexander Marrow, Darya Korsunskaya and Polina Devitt

November 25, 2022 7:31 AM PST · Updated 8 months ago





Europe

Putin, Kudrin touch on future of Yandex in late-night meeting -sources

Reuters

November 25, 2022 4:19 AM PST · Updated 8 months ago



[1/2] The logo of Russian internet group Yandex is pictured at the company's headquarter in Moscow, Russia October 4, 2018. REUTERS/Shamil Zhumatov



Kremlin Ally Kudrin Confirms Move to Tech Giant Yandex

Dec. 5, 2022





YANDEX SERVICES SOURCE CODE LEAK

SHORT OVERVIEW OF BREACH CONTENTS

PUBLISHED THU, JAN 26, 2023 BY ARSENIY SHESTAKOV

Just a few hours ago I found [mention on Twitter](#) that proprietary source code of Russian giant Yandex been leaked on online community called *BreachForums*. In this post I'll share results of my **friend** digging into said archives.

Important details about torrent:

- It just content of repository without anything else.
- All files are dated back to [24 February 2022](#).
- It does not contain git history, mostly just code
- No pre-built binaries for most of software with only few exceptions
- There are no pre-trained ML models with some exceptions

Russian Billionaires Line Up to Buy Yandex – Reports

May 4, 2023



MOST READ JUST IN

- OVERNIGHT STRIKE**
Russia Says Ukrainian Drones Hit Moscow, Crimea
- NO PASSAGE**
Russia Blocks Cargo Ship Over 'Traces'
- POLITICAL PRISONER**
Navalny Ally Jailed 9 Years for
- MORE MANPOWER**
Russia Raises Upper-Age Limit for Reservists
- MONEY DRAIN**



Technology

Russia's Yandex fined for refusing to share user information with security services

Reuters

June 18, 2023 3:20 PM PDT · Updated a month ago

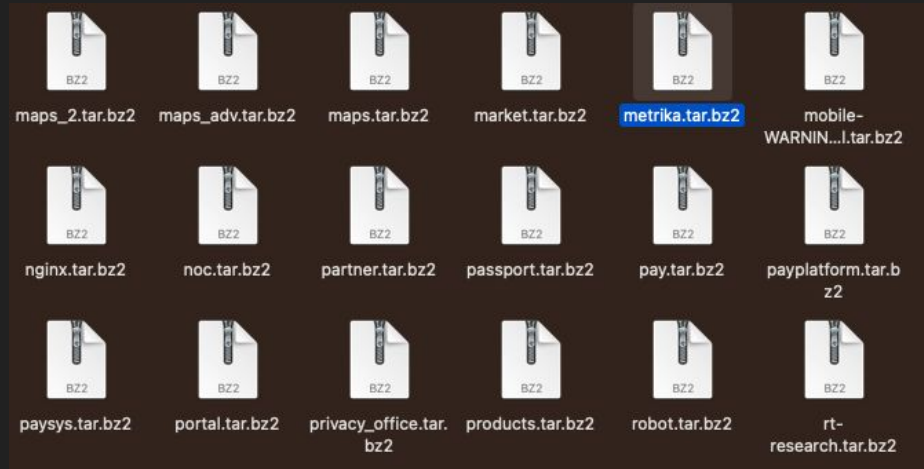


Roadmap

- Background on Yandex Leak
- Dive into code:
 - What data Yandex is collecting
 - What Yandex is doing with that data
 - Who Yandex is sharing that data with
- Conclusions and wrap up
- Q&A

Yandex Codebase

Codebase



Metrika

Metrika



Solutions

Features

Verticals

Resources

Supercharge app metrics with data insights

with a one-stop solution for analytics and marketing

Get started

Take a tour

Yandex Metrika

Features

Resources

Pricing

All-Round Web Analytics

From traffic trends to mouse movements – get a comprehensive understanding of your online audience and drive business growth.

Get started

Try live demo

Example Raw Data Fields that AppMetrica Logs

```
analytics > appmetrica-location-log-anonymizer > ☰ convert_log.yml
89 insert into `//home/metrika-logs/anonym-appmetrica-location-log/1d/{table_date}`
90 with truncate
91 select
92     String::HexEncode(Digest::Blake2B(`DeviceID`, seed)) as `DeviceID`,
93     String::HexEncode(Digest::Blake2B(`ADVID`, seed)) as `ADVID`,
94     String::HexEncode(Digest::Blake2B(`IFA`, seed)) as `IFA`,
95     String::HexEncode(Digest::Blake2B(`UUID`, seed)) as `UUID`,
96     String::HexEncode(Digest::Blake2B(`AndroidID`, seed)) as `AndroidID`,
97     `APIKey`,
98     `AppBuildNumber`,
99     `AppFramework`,
100    `AppID`,
101    `AppPlatform`,
102    `AppVersionName`,
103    `Cells_AreConnected`,
104    `Cells_CellsIDs`,
105    `Cells_CountriesCodes`,
106    `Cells_Lacs`,
107    `Cells_LastVisibleTimeOffset`,
108    `Cells_OperatorsIDs`,
109    `Cells_OperatorsNames`,
110    `Cells_PhysicalsCellsIDs`,
111    `Cells_SignalsStrengths`,
112    `Cells_Types`,
113    `ChargeType`,
114    `ClientIP`,
115    `ClientIPHash`,
116    `CollectTimestamp`,
117    `CollectTimestampBootOffset`, |
118    `CollectionMode`,
119    `DeviceType`,
120    `EventID`,
121    `IncrementalID`,
122    `IsExtraLocationEvent`,
123    `IsRooted`,
124    `KitBuildNumber`,
125    `KitBuildType`,
126    `KitVersion`,
```

Ln 117

```
analytics > appmetrica-location-log-anonymizer > ☰ convert_log.yml
122     `IsExtraLocationEvent`,
123     `IsRooted`,
124     `KitBuildNumber`,
125     `KitBuildType`,
126     `KitVersion`,
127     `Latitude`,
128     `LatitudeLBS`,
129     `LocationAltitude`,
130     `LocationDirection`,
131     `LocationEnabled`,
132     `LocationPrecision`,
133     `LocationPrecisionLBS`,
134     `LocationSource`,
135     `LocationSpeed`,
136     `LocationTimestamp`,
137     `LocationTimestampBootOffset`,
138     `Longitude`,
139     `LongitudeLBS`,
140     `OSApiLevel`,
141     `OSVersion`,
142     `OperatingSystem`,
143     `OriginalCollectTimestamp`,
144     `OriginalLocationTimestamp`,
145     `ReceiveTimestamp`,
146     `RequestID`,
147     `SendTimestamp`,
148     `Wifi_AreConnected`,
149     `Wifi_LastVisibleTimeOffset`,
150     `Wifi_Macs`,
151     `Wifi_SignalsStrengths`,
152     `Wifi_Ssids`,
153     `_logfeller_index_bucket`,
154     `_logfeller_timestamp`,
155     `_rest`,
156     `_stbx`,
157     `iso_eventtime`,
158     `source_uri`,
159     `subkey`,
160     `-----`
```

Anonymized identifiers

```
91      select
92          String::HexEncode(Digest::Blake2B(`DeviceID`, seed)) as `DeviceID`,
93          String::HexEncode(Digest::Blake2B(`ADVID`, seed)) as `ADVID`,
94          String::HexEncode(Digest::Blake2B(`IFA`, seed)) as `IFA`,
95          String::HexEncode(Digest::Blake2B(`UUID`, seed)) as `UUID`,
96          String::HexEncode(Digest::Blake2B(`AndroidID`, seed)) as `AndroidID`,
```

Location Fields

```
126 `KitVersion`,  
127 `Latitude`,  
128 `LatitudeLBS`,  
129 `LocationAltitude`,  
130 `LocationDirection`,  
131 `LocationEnabled`,  
132 `LocationPrecision`,  
133 `LocationPrecisionLBS`,  
134 `LocationSource`,  
135 `LocationSpeed`,  
136 `LocationTimestamp`,  
137 `LocationTimestampBootOffset`,  
138 `Longitude`,  
139 `LongitudeLBS`,  
140 `OSApiLevel`,  
141 `OSVersion`,  
142 `OperatingSystem`,  
143 `OriginalCollectTimestamp`,  
144 `OriginalLocationTimestamp`.
```

Wifi Fields Collected By AppMetrica

147	<code>`SendTimestamp`,</code>
148	<code>`Wifi_AreConnected`,</code>
149	<code>`Wifi_LastVisibleTimeOffset`,</code>
150	<code>`Wifi_Macs`,</code>
151	<code>`Wifi_SignalsStrengths`,</code>
152	<code>`Wifi_Ssids`,</code>
153	<code>` logfeller index bucket`,</code>

Those fields in Crypta

```
graph > fuzzy > lib > yql > ≡ export_ssid_devid_day_table.yql
```

```
26
27  $list_metrika_log = (
28      select coalesce(DeviceID, "") as DeviceID,
29             coalesce(OriginalDeviceID, "") as OriginalDeviceID,
30             $MakeStringList(Wifi_Macs) as Wifi_Macs,
31             $MakeStringList(Wifi_Ssids) as Wifi_Ssids,
32             $MakeIntList(Wifi_SignalsStrengths) as Wifi_SignalsStrengths,
33             $MakeIntList(Wifi_AreConnected) as Wifi_AreConnected
34      from `{source_mmetric_table}`
35      where DeviceID is not null
36  );
```


Dev Id and SSID Associated with Yandex UID

```
graph > fuzzy > lib > yql > ≡ export_ssid_yuids.yql
```

```
6  $mobile_all_table = (  
7  |   select distinct mmetric_devid, ssid  
8  |   from concat({sources})  
9  | );  
10  
11 $mmetric_to_devid = (  
12 |   select mmetric_devid, devid,  
13 |   |   coalesce(cast(yuid as uint64), 0) as yuid  
14 |   from `{source_nolimit_table}`  
15 | );  
16
```

Click Event Data Being Matched to Existing Users

```
core > programs > clicklogd-mobile > src > C event_indexed_pool.h > TEventIndexedPool > GetIndex<TMatchCriteria>()
58 private:
59     template <class TMatchCriteria>
60     TIndex<TMatchCriteria>& GetIndex() {
61         if constexpr (std::is_same_v<TMatchCriteria, NMatchCriteria::TAndroidId>) {
62             return AndroidId_;
63         } else if constexpr (std::is_same_v<TMatchCriteria, NMatchCriteria::TAndroidIdMd5>) {
64             return AndroidIdMd5_;
65         } else if constexpr (std::is_same_v<TMatchCriteria, NMatchCriteria::TAndroidIdSha1>) {
66             return AndroidIdSha1_;
67         } else if constexpr (std::is_same_v<TMatchCriteria, NMatchCriteria::TDeviceIdHash>) {
68             return DeviceIdHash_;
69         } else if constexpr (std::is_same_v<TMatchCriteria, NMatchCriteria::TFingerprint>) {
70             return Fingerprint_;
71         } else if constexpr (std::is_same_v<TMatchCriteria, NMatchCriteria::TGoogleAid>) {
72             return GoogleAid_;
73         } else if constexpr (std::is_same_v<TMatchCriteria, NMatchCriteria::TGoogleAidMd5>) {
74             return GoogleAidMd5_;
75         } else if constexpr (std::is_same_v<TMatchCriteria, NMatchCriteria::TGoogleAidSha1>) {
76             return GoogleAidSha1_;
77         } else if constexpr (std::is_same_v<TMatchCriteria, NMatchCriteria::TIfa>) {
78             return Ifa_;
79         } else if constexpr (std::is_same_v<TMatchCriteria, NMatchCriteria::TIfaMd5>) {
80             return IfaMd5_;
81         } else if constexpr (std::is_same_v<TMatchCriteria, NMatchCriteria::TIfaSha1>) {
82             return IfaSha1_;
83         } else if constexpr (std::is_same_v<TMatchCriteria, NMatchCriteria::TWindowsAid>) {
84             return WindowsAid_;
85         } else if constexpr (std::is_same_v<TMatchCriteria, NMatchCriteria::TWindowsAidMd5>) {
86             return WindowsAidMd5_;
87         } else if constexpr (std::is_same_v<TMatchCriteria, NMatchCriteria::TWindowsAidSha1>) {
88             return WindowsAidSha1_;
89         } else if constexpr (std::is_same_v<TMatchCriteria, NMatchCriteria::TYmTrackingId>) {
90             return YmTrackingId_;
91         }
92     }
93 }
```

Socio-Demographic Attributes for DevID being Updated

```
core > programs > socdem-updaterd-mobile > src > C: UserIdAndInfoParser.cpp > ...
45     {"0_17", AgeIntervalsCrypta::LessThan18},
46     {"18_24", AgeIntervalsCrypta::Between18and24},
47     {"25_34", AgeIntervalsCrypta::Between25and34},
48     {"35_44", AgeIntervalsCrypta::Between35and44},
49     {"45_54", AgeIntervalsCrypta::Between45and54},
50     {"55_99", AgeIntervalsCrypta::MoreThan55}
51 };
52     ::setValue(value, exact_socdem_node, key, json_keys_to_ages_intervals);
53 }
54
55 void UserIdAndInfoParser::setValue(
56     SexTypesCrypta & value,
57     const NYT::TNode & exact_socdem_node,
58     const TString key)
59 {
60     static const std::map<TString, SexTypesCrypta> json_keys_to_sex_types =
61     {
62         {"f", SexTypesCrypta::Female},
63         {"m", SexTypesCrypta::Male}
64     };
65     ::setValue(value, exact_socdem_node, key, json_keys_to_sex_types);
66 }
67
68 std::string UserIdAndInfoParser::parse(const NYT::TNode & user_record)
69 {
70     const TString & device_id = user_record["appmetrica_devid"].AsString();
71     UserInfo user_info;
72
73     const auto & exact_socdem = user_record["exact_socdem"];
74     setValue(user_info.age, exact_socdem, "age_segment");
75     setValue(user_info.sex, exact_socdem, "gender");
76
77     static const auto tail = getConstTail();
78
79     std::ostringstream buffer;
80     buffer <<
81         sipHash64(device_id.data(), device_id.size()) << '\t' <<
82         static_cast<int>(user_info.age) << '\t' <<
```



Create segments based on offline and online data

Create Segment

Лояльные клиенты

Охват: **582 685**

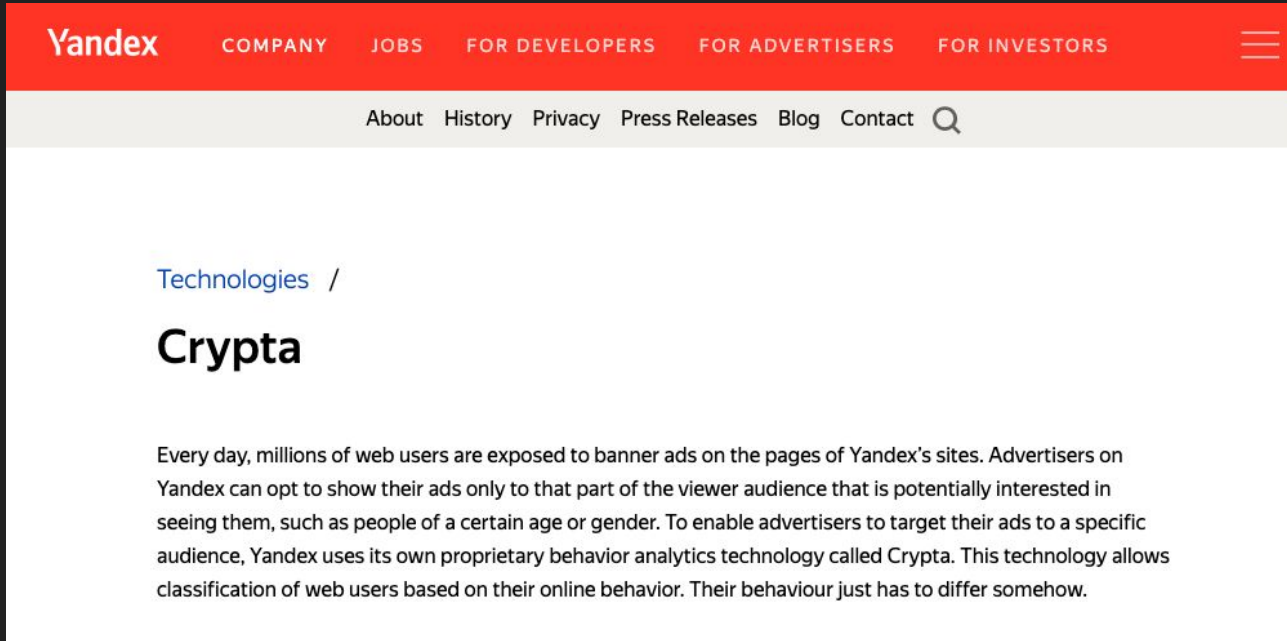
Возраст: 5% (18-24), 18% (25-34), 35% (35-44), 29% (45-54), 13% (55-64)

Женщины: **73%** | Мужчины: **27%**

Города	Устройства
Москва: 38%	ПК: 73%
Санкт-Петербург: 30%	Смартфоны: 17%
Воронеж: 26%	Планшеты: 9%
Екатеринбург: 3%	ТВ: 1%
Ростов-на-Дону: 2%	
Краснодар: 1%	


Crypta

Crypta



The screenshot shows the Yandex website's header with navigation links: COMPANY, JOBS, FOR DEVELOPERS, FOR ADVERTISERS, and FOR INVESTORS. Below the header is a secondary navigation bar with links for About, History, Privacy, Press Releases, Blog, and Contact, along with a search icon. The main content area features a breadcrumb trail 'Technologies /' followed by the title 'Crypta'. The text below explains that Yandex uses its proprietary behavior analytics technology, Crypta, to target ads to specific audiences based on their online behavior.

Yandex COMPANY JOBS FOR DEVELOPERS FOR ADVERTISERS FOR INVESTORS

About History Privacy Press Releases Blog Contact 

Technologies /

Crypta

Every day, millions of web users are exposed to banner ads on the pages of Yandex's sites. Advertisers on Yandex can opt to show their ads only to that part of the viewer audience that is potentially interested in seeing them, such as people of a certain age or gender. To enable advertisers to target their ads to a specific audience, Yandex uses its own proprietary behavior analytics technology called Crypta. This technology allows classification of web users based on their online behavior. Their behaviour just has to differ somehow.

Example Segments

CRYPTA	
> test	
age_segment_18_20.py	
alice_users.py	
apartment_room_number.py	
apps_users.py	
artists.py	
auto_interactions.py	
avia_travellers.py	
bank_cards.py	
bought_two_tickets.py	
business_travellers.py	
children_age_segment_clarification.py	
compulsory_auto_insurance.py	
connection_type.py	
console_gamers.py	
contest.py	
devices_without_google_services.py	
digital_viewers.py	
direct_clients_by_industry.py	
direct_product_users.py	
disk_users.py	
ecommerce_owners.py	
edadeal_offline_purchases_lal.py	
expensive_car_customers.py	
film_lovers_by_genres.py	
gas_stations.py	
industry_representatives.py	
kbt_customers.py	
kfc_visitors.py	
kinopoisk_logins.py	
kinopoisk_movie_watchers.py	

CRYPTA	
kinopoisk_movie_watchers.py	
kz_users.py	
laptop_users.py	
logged_in_for_plus.py	
longterm_interest_mobile_gamers.py	
loyal_to_launcher_install.py	
macos_users.py	
mail_data.py	
manufacturer_phone_owners.py	
marketplaces_tv_users.py	
mobile_gamers.py	
mobile_operators_users_by_prefix.py	
mobile_operators_users.py	
multidevice_puid.py	
multidevice.py	
music_genres_listeners.py	
nestle_regions.py	
phone_buyers.py	
phone_owners.py	
phone_with_esim_owners.py	
potential_aon_android_users.py	
potential_aon_ios_users.py	
preinstalled_apps.py	3
prism.py	
proleads.py	
realty_interactions.py	
recent_passport_accounts.py	
score_users_for_telephony.py	
searched_for_phone_numbers.py	
searched_radisson_on_maps.py	

CRYPTA	
mobile_operators_users_by_prefix.py	
mobile_operators_users.py	
multidevice_puid.py	
multidevice.py	
music_genres_listeners.py	
nestle_regions.py	
phone_buyers.py	
phone_owners.py	
phone_with_esim_owners.py	
potential_aon_android_users.py	
potential_aon_ios_users.py	
preinstalled_apps.py	
prism.py	
proleads.py	
realty_interactions.py	
recent_passport_accounts.py	
score_users_for_telephony.py	
searched_for_phone_numbers.py	5
searched_radisson_on_maps.py	
seo_specialists.py	
seo_users.py	
smart_gadgets_customers.py	
smokers.py	
summer_residents.py	
travellers.py	
video_bloggers.py	
want_to_change_the_provider.py	
webmaster.py	
widgets.py	
with_children_by_ages.py	

Example Segments

- 🔗 [smart_gadgets_customers.py](#)
- 🔗 [smokers.py](#)
- 🔗 [summer_residents.py](#)
- 🔗 [travellers.py](#)
- 🔗 [video_bloggers.py](#)
- 🔗 [want_to_change_the_provider.py](#)
- 🔗 [webmaster.py](#)
- 🔗 [widgets.py](#)
- 🔗 [with_children_by_ages.py](#)
- ☰ [va_make](#)

Travellers

profile > runners > segments > lib > coded_segments >  travellers.py > ...

```
82
83 INSERT INTO `{output_table}` WITH TRUNCATE
84 SELECT
85     id,
86     id_type,
87     segment_name
88 FROM(
89     SELECT
90         crypta_id AS id,
91         'crypta_id' AS id_type,
92         CASE
93             WHEN Geo::RoundRegionById(region, "country").id != Geo::RoundRegionById(CAST(main_region AS Int32), "country").id THEN 'international'
94             ELSE 'domestic'
95         END AS segment_name,
96         MAX(`date`) AS last_seen,
97         MIN(`date`) AS first_seen,
98         region,
99         week_end_date,
100     FROM $travell_visits
101     GROUP BY region, main_region, crypta_id, week_end_date
102 )
103 WHERE
104     last_seen <= week_end_date AND
105     DateTime::ToDays(DateTime::MakeTimestamp($parse(last_seen)) - DateTime::MakeTimestamp($parse(first_seen))) > 0 AND
106     DateTime::ToDays(DateTime::MakeTimestamp($parse(week_end_date)) - DateTime::MakeTimestamp($parse(first_seen))) <= 7
107 GROUP BY id, id_type, segment_name
108 """"
109
110
```

Mail Data

```
profile > runners > segments > lib > coded_segments > mail_data.py > ...
12
13
14 segment_query = """
15 INSERT INTO `{output_table}` WITH TRUNCATE
16 SELECT id, id_type, segment_name
17 FROM `{mail_data_table}`;
18
19 INSERT INTO `{sample_table}` WITH TRUNCATE
20 SELECT
21     yandexuid,
22     segment_name
23 FROM (
24     SELECT matching.yandexuid AS yandexuid, mail_data.segment_name AS segment_name
25     FROM `{mail_data_table}` AS mail_data
26     INNER JOIN `{indevice_yandexuid_matching}` AS matching
27     USING (id, id_type)
28 )
29 GROUP BY yandexuid, segment_name
30 """
31
32
33 class PrepareMailSampleForLalSegments(RegularSegmentBuilder):
34     keyword = 549
35     name_segment_dict = {
36         'aviaticket': 1404,
37         'boardingpass': 1405,
38         'hotel': 1406,
39     }
```

Gas Stations

```
profile > runners > segments > lib > coded_segments > gas_stations.py > ...
__
92 class ProcessedDeepVisitLogForGasStations(DayProcessor):
93     def requires(self):
94         return deep_visits.org_visits_deep_external_input(self.date)
95
96     def process_day(self, inputs, output_path):
97
98         self.yql.query(
99             gas_stations_query_template.format(
100                 organization_categories=config.ORGANIZATION_CATEGORIES,
101                 deep_visits=inputs.table,
102                 matching_idfa=get_matching_table('idfa', 'crypta_id'),
103                 matching_gaid=get_matching_table('gaid', 'crypta_id'),
104                 name_to_variable=',\n'.join(
105                     [u'("{} ", "{}")'.format(key, value)
106                     for key, value in name_to_variable.iteritems()]
107                 ),
108                 output_table=output_path,
109             ),
110             transaction=self.transaction,
111         )
112
```

Example ML Model Types

```
▼ CRYPTA
  legal_entities_model_application.py
  legal_entities_model_training.py
  legal_office_visits_model_application.py
  legal_office_visits_model_training.py
  market_model_application.py
  market_model_training.py
  market_rfm_model_application.py
  market_rfm_model_training.py
  marriage_model_application.py
  marriage_model_training.py
  medical_clinic_model_application.py
  medical_clinic_model_training.py
  mortgage_approval_model_application.py
  mortgage_approval_model_training.py
  online_cinema_model_application.py
  online_cinema_model_training.py
  online_payment_model_application.py
  online_payment_model_training.py
  online_sales_register_model_application.py
  online_sales_register_model_training.py
  online_shopping_model_application.py
  online_shopping_model_training.py
  pharmacy_model_application.py
  pharmacy_model_training.py
  realty_visit_model_application.py
  realty_visit_model_training.py
  tv_viewers_model_application.py
  tv_viewers_model_training.py
  windows_installation_model_application.py
  windows_installation_model_training.py
  va make
```

Basic example of household details

```
graph > metrics > household > query.sql
44     FROM $composition
45     GROUP BY $size_to_range(size) AS key;
46 END DEFINE;
47
48 DEFINE SUBQUERY $hh_size_by_crypta_id($title, $predicat) AS
49     SELECT
50         ($title || key) AS key,
51         COUNT(1) AS hh_size_by_crypta_id
52     FROM $composition
53     WHERE $predicat(size, socdems)
54     GROUP BY CAST(Yson::GetLength(Yson::Lookup(data, 'crypta_ids')) AS String) AS key;
55 END DEFINE;
56
57 DEFINE SUBQUERY $hh_by_socdems($title, $predicat) AS
58     $hh_socdem = (
59         SELECT
60             hhid,
61             size,
62             IF((Yson::LookupInt64(info, 'female') != 0), 'female', Null) AS has_female,
63             IF((Yson::LookupInt64(info, 'male') != 0), 'male', Null) AS has_male,
64             IF((Yson::LookupInt64(info, 'grand') != 0), 'grand', Null) AS has_old,
65             IF((Yson::LookupInt64(info, 'child') != 0), 'child', Null) AS has_child
66         FROM $composition
67         WHERE $predicat(size, socdems)
68     );
69
70     SELECT ($title || groups) AS key, hh_c AS hh_socdem_count
71     FROM (
72         SELECT groups, SUM(size) AS hh_c
73         FROM $hh_socdem
74         GROUP BY String::JoinFromList(
75             listSort(AsList(has female, has male, has old, has child)),
76             '_' ) AS groups
77     ) WHERE groups != "";
78 END DEFINE;
79
```

AppMetrica being used to pull wifi connection types:

```
profile > runners > segments > lib > coded_segments > connection_type.py > ...
59 connection_type_query = """
60 INSERT INTO `{output_table}` WITH TRUNCATE
61 SELECT
62     id,
63     'mm_device_id' AS id_type,
64     CASE
65         WHEN types == AsSet('3g') THEN '3g'
66         WHEN types == AsSet('4g') THEN '4g'
67         ELSE '3g_4g'
68     END AS segment_name
69 FROM (
70     SELECT
71         id,
72         ToSet(AGGREGATE_LIST_DISTINCT(segment_name)) AS types
73     FROM `{input_table}`
74     GROUP BY id
75 );
76 """
77
78
79 class ConnectionType(RegularSegmentBuilder):
80     name_segment_dict = {
81         '3g': (557, 17823841),
82         '4g': (557, 17823853),
83         '3g_4g': (557, 17823847),
84     }
85
86     number_of_days = 35
87
88     def requires(self):
89         return {
90             'AppMetrica': LogProcessor(
91                 ProcessAppMetricaForConnectionType,
92                 self.date,
93                 self.number_of_days,
94             ),
95         }
```

AppMetrica data being used to separate users with common SSIDs (wifi networks)

```
class ImportSsidMobileMetrikaTask(BaseTask):
    date = DateParameter()
    SSID_THRESHOLD = 20
    YUID_THRESHOLD = 20
    DAYS_IN_MONTH = 7

    def requires(self):
        """
        this tasks must be done to complete this task
        """
        task_list = [
            ImportSsidMobileMetrikaDayTask(date=self.date, target_date=target_date, ssid_threshold=self.SSID_THRESHOLD)
            for target_date in days_range_back(self.date, self.DAYS_IN_MONTH)
        ]
        return task_list
```


AppMetrica data being used to separate users with common SSIDs (wifi networks)

```
def _run(self):
    self.yt.create_table_with_schema(
        self.destination, self.destination_schema, strict=True, recreate_if_exists=True
    )
    with self.yt.TempTable() as unexploded, self.yt.TempTable() as not_unique:
        self.yql.execute(self.query(unexploded), syntax_version=1)
        run_native_reduce(
            reducer_name="NCommonWifiAP::TExploder",
            source=unexploded,
            destination=not_unique,
            proxy=self.yt.proxy,
            transaction=self.yt.transaction_id,
            pool=conf.Yt.POOL,
            title="Explode yandexuids with common wifi access point",
            reduce_by=["ssid"],
        )
        yuid_pair = [conf.Constants.YUID_LEFT, conf.Constants.YUID_RIGHT]
        self.yt.run_sort(not_unique, not_unique, sort_by=yuid_pair)
        run_native_reduce(
            reducer_name="NCommonWifiAP::TUnique",
            source=not_unique,
            destination=self.destination,
            proxy=self.yt.proxy,
            transaction=self.yt.transaction_id,
            pool=conf.Yt.POOL,
            title="Make yandexuids with common wifi access point unique",
            reduce_by=yuid_pair,
        )
        self.yt.run_sort(self.destination, sort_by=yuid_pair)
    self.yt.set(self.destination + "/@generate_date", self.date.isoformat())
```


Sources

```
graph > fuzzy > lib > config.py > GeoPaths
```

```
59
```

```
60 class SourceTypes(object):
```

```
61     EMAIL_LOGIN = "EMAIL_LOGINS"
```

```
62     EMAIL_SIMILAR = "EMAIL_SIMILAR"
```

```
63     GEO_HOMEWORK = "GEO_HOMEWORK"
```

```
64     HOUSEHOLD = "HOUSEHOLD"
```

```
65     REQANS_LOG = "REQANS_LOG" ← Search Data
```

```
66     SSID = "SSID" ← Wifi
```

```
67
```

```
68
```

Yandex IDs Associated with Email

```
class EmailPaths(object):
    ROOT = ROOT
    # Emails
    BASE = "{root}/email".format(root=ROOT)

    ALL_EMAILS_TABLE = "{base}/all_emails".format(base=BASE)
    ALL_EMAIL_LOGINS_TABLE = "{base}/all_email_logins".format(base=BASE)
    ALL_EMAILS_SORTED_BY_LOGIN = "{base}/all_email_logins.sorted_by_login".format(base=BASE)
    ALL_EMAIL_LOGINS_PAIRS_TABLE = "{base}/all_email_logins.pairs".format(base=BASE)
    ALL_EMAILS_GROUPED_BY_LOGIN = "{base}/all_email_logins.groups".format(base=BASE)
    ALL_YUID_PAIRS_FROM_EMAIL_LOGIN = "{base}/all_yuid_pairs_from_email_logins_matching".format(base=BASE)
    ALL_YUID_PAIRS_FROM_SIMILAR_EMAILS = "{base}/all_yuid_pairs_from_similar_emails".format(base=BASE)

    ALL_EMAILS_TABLE_SCHEMA = {"email": "string", "yuids": "any"}
    ALL_EMAIL_LOGINS_TABLE_SCHEMA = {"login": "string", "email": "string", "yuids": "any"}
    ALL_EMAILS_SORTED_BY_LOGIN_SCHEMA = {"login": "string", "email": "string", "yuids": "any"}
    ALL_EMAILS_GROUPED_BY_LOGIN_SCHEMA = {"login": "string", "all_emails": "any", "howmany": "uint64"}
    ALL_EMAIL_LOGINS_PAIRS_TABLE_SCHEMA = {
        "email_1": "string",
        "email_2": "string",
        "login": "string",
        "yuids_1": "any",
        "yuids_2": "any",
    }

    ALL_YUID_PAIRS_FROM_EMAIL_LOGIN_SCHEMA = {
        Constants.YUID_LEFT: ("uint64", True),
        Constants.YUID_RIGHT: ("uint64", True),
        "match": "any",
    }

    ALL_YUID_PAIRS_FROM_SIMILAR_EMAILS_SCHEMA = {
        Constants.YUID_LEFT: "uint64",
        Constants.YUID_RIGHT: "uint64",
        "email_left": "string",
        "email_right": "string",
        "fragment": "string",
    }
}
```

Login Data

graph > fuzzy > lib > tasks > sources > visitlog_logins > extract.py > ...

```
59     def filter_rare_logins_options(self):
60         return TFilterRareLoginsOptions(Threshold=self.threshold).SerializeToString()
61
62     @property
63     def filter_keys_options(self):
64         return TFilterKeysOptions(
65             Keywords=[
66                 "login",
67                 "user",
68                 "userid",
69                 "clientid",
70                 "uid",
71                 "email",
72                 "emailhash",
73                 "\u043b\u043e\u0433\u0438\u043d\u043d",
74                 "computerid",
75                 "cid",
76                 "suserid",
77             ]
78         ).SerializeToString()
79
80
```

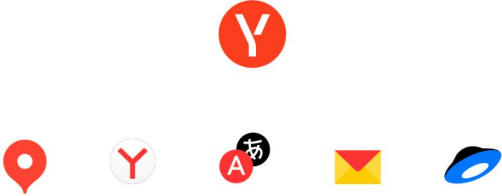
Extracting
multiple types
of identifiers

Passport

passport.yandex.com/registration?retpath=https://audience.yandex.com/

CyberChef confiant-inc/priva... Tracker DBs Reading Useful snippets https://cmplist.co... https://vendor-list... TCF Consent String De... Technique protect...

Already have an a



With a single account, you can search, write emails, save and share your files, find stuff you want, get directions and use other services on all your devices and platforms

Registration

First name

Last name

Enter a login

Enter a password

Confirm password

Mobile phone number

Register

[Help](#) © 2023, Yandex

Passport User ID Associated with Phone

```
graph > data_import > passport > lib > query > passport.sql
33
34 $out_login_tbl = $soup_output_dir || $edge(IdType::PUID(), IdType::PHONE(), SourceType::PASSPORT_PROFILE(), LogSource::PASSPORT_PHONE_DUMP());
35 INSERT INTO $out_login_tbl WITH TRUNCATE
36 SELECT
37     id1,
38     IdType::PUID() AS id1Type,
39     id2,
40     IdType::PHONE() AS id2Type,
41     SourceType::PASSPORT_PROFILE() AS sourceType,
42     LogSource::PASSPORT_PHONE_DUMP() AS logSource,
43     ListCreate(String) AS dates
44 FROM (
45     SELECT DISTINCT puid, phone
46     FROM $log FLATTEN LIST BY phones AS phone
47 ) WHERE Identifiers::IsSignificantPhone(phone)
48 GROUP BY
49     puid AS id1,
50     phone AS id2
51 ;
52
```

```
graph > data_import > passport > tests > fixtures > {} passport.json > ...
```

```
1  {"uid": "11111", "login": "aashinova"}
2  {"uid": "11112", "login": "andrei-ponomareff-1997"}
3  {"uid": "11113", "login": "anoshko-av"}
4  {"uid": "11114", "login": "bars12@161.ru"}
5  {"uid": "11115", "login": "ev0ngertlt"}
6  {"uid": "11116", "login": "evarcher"}
7  {"uid": "11117", "login": "lagutin2008"}
8  {"uid": "11118", "login": "login-for-avito"}
9  {"uid": "11119", "login": "modsever"}
10 {"uid": "11120", "login": "mouradian"}
11 {"uid": "11121", "login": "perschina-olga2013"}
12 {"uid": "11122", "login": "r.amiraslanov@dveri.ru"}
13 {"uid": "11123", "login": "saprovec2015"}
14 {"uid": "11124", "login": "stoltat"}
15 {"uid": "11125", "login": "sveta-aleshina2015"}
16 {"uid": "11126", "login": "watchradius"}
17 {"uid": "11127", "login": "watchradius"}
18 {"uid": "11128", "login": "sveta-aleshina2015"}
19 {"uid": "11129", "login": "watchradiusmob"}
20 {"uid": "123456", "login": "abc123", "phone_numbers": []}
21 {"uid": "123457", "login": "abc127", "phone_numbers": ["+1234567890", "+71111234567"]}
22 {"uid": "123458", "login": "", "phone_numbers": ["+9393939393", "+712020201"]}
23 {"uid": "134614616", "login": "roscosh8"}
24 {"uid": "134648582", "login": "e222mn"}
25 {"uid": "15033290", "login": "mouradian"}
26 {"uid": "194502233", "login": "ingvr80"}
27 {"uid": "2687", "login": "govshit"}
28 {"uid": "2687", "login": "GOVSHIT"}
29 {"uid": "76667777", "login": "g8jkqqaah"}
30 {"uid": "766679666", "login": "d6fqqaah"}
31 {"uid": "766679777", "login": "mdmozn45"}
32
```

Crypta - Geo graphs

Using lat/long data associated with “predicted home”, linked to Yandex UID

```
graph > fuzzy > lib > tasks > sources > geo > C geo_operations.h
59 void Do(TTableReader<TNode>* input, TTableWriter<TGeoSquare>* output)
60 override {
61     for (; input->IsValid(); input->Next()) {
62         const auto& row = input->GetRow();
63         if (not IsValidRow(row)) {
64             continue;
65         }
66
67         const ui64 yandexuid = FromString<ui64>(row["yandexuid"].AsString());
68         const auto& homeCoordinates = row["predicted_home"];
69         const auto latitude = homeCoordinates["latitude"].AsDouble();
70         const auto longitude = homeCoordinates["longitude"].AsDouble();
71         const auto& square = computeSquare({.Lat = latitude, .Lon = longitude}, State->radius());
72
73         /**
74          * +- -
75          * +- -
76          * + - -
77          */
78         for (int beltOffset : {-1, 0, 1}) {
79             for (int sqOffset : {-1, 0}) {
80                 if (beltOffset == -1 && sqOffset == 0) {
81                     continue;
82                 }
83                 const ui64 square_idx = ConvertSquareToIdx({.Belt = square.Belt + beltOffset, .Sq = square.Sq + sqOffset});
84                 TGeoSquare out;
85                 out.set_yandexuid(yandexuid);
86                 out.set_lat(latitude);
87                 out.set_lon(longitude);
88                 out.set_squareidx(square_idx);
89                 output->AddRow(out);
90             }
91         }
92     }
93 }
```


Crypta - Geo graphs

Then using that data to find literal neighbors within a certain radius of that home

```
graph > fuzzy > lib > tasks > sources > geo > C geo_operations.h
106
107 class TFindNeighbors: public IReducer<TTableReader<TGeoSquare>, TTableWriter<TNeighborsDistance>> {
108 public:
109     TFindNeighbors()
110         : State()
111     {
112     }
113
114     TFindNeighbors(const TBuffer& buffer)
115         : State(buffer)
116     {
117     }
118
119     void Do(TTableReader<TGeoSquare>* input, TTableWriter<TNeighborsDistance>* output) override {
120         const double radius = State->radius();
121         TVector<TGeoSquare> candidates;
122         for (; input->IsValid(); input->Next()) {
123             const auto& row = input->GetRow();
124             candidates.push_back(row);
125         }
126         for (auto i : xrange(candidates.size())) {
127             for (auto j : xrange(i + 1, candidates.size())) {
128                 const auto& left = candidates.at(i);
129                 const auto& right = candidates.at(j);
130                 if (left.yandexuid() == right.yandexuid()) {
131                     continue;
132                 }
133                 double distance = computeDistance({.Lat = left.lat(), .Lon = left.lon()}, {.Lat = right.lat(), .Lon = right.lon()});
134                 if (distance > radius) {
135                     continue;
136                 }
137                 TNeighborsDistance out;
138                 out.set_distance(distance);
139                 out.set_yandexuidleft(Min(left.yandexuid(), right.yandexuid()));
140                 out.set_yandexuidright(Max(left.yandexuid(), right.yandexuid()));
141                 output->AddRow(out);
142             }
143         }
144     }
145 }
```

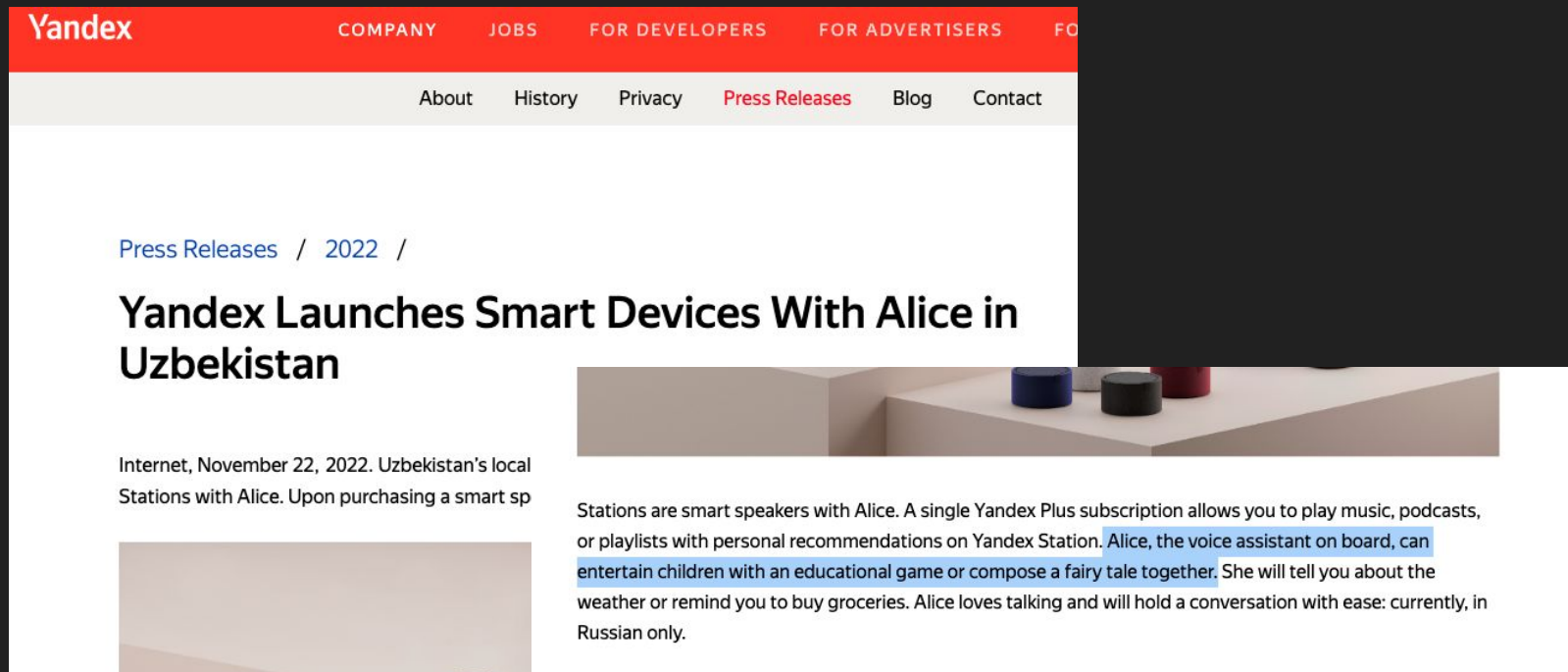

AppMetrica and Taxi data being used generate segments about households with children:

```
self.yql.query(  
    app_metrica_query.format(  
        devid_by_app_table=self.input()['DevidByApp'].table,  
        output_table=with_children_by_app_table,  
        app_to_segment_name='\n'.join(app_segment_name_tuples),  
    ),  
    transaction=self.transaction,  
)  
  
def build_segment(self, inputs, output_path):  
    with self.yt.TempTable() as taxi_puid_table, \  
        self.yt.TempTable() as app_metrica_table:  
        self.yt.run_map(  
            extract children from taxi,  
            inputs['TaxiData'].table,  
            taxi_puid_table,  
        )  
  
        self.prepare_with_children_by_app(app_metrica_table)  
  
        self.yql.query(  
            with_children_query_template.format(  
                metrics_table=inputs['ProcessedMetrics'].table,  
                reqans_table=inputs['ProcessedReqans'].table,  
                app_metrica_table=app_metrica_table,  
                taxi_data_table=taxi_puid_table,  
                id_to_crypta_id_table=config.VERTICALS_NO_MULTI_PROFILE,  
                crypta_id_to_hhid_table=config.HOUSEHOLD_CRYPTA_ID_TO_HHID,  
                yandexuid_to_hhid_table=config.HOUSEHOLD_REVERSED_TABLE,  
                hhid_to_yandexuid_table=config.HOUSEHOLD_ENRICH_TABLE,  
                output_table=output_path,  
            ),
```

ID mapping associations:

```
taxi_data_table=taxi_puid_table,  
id_to_crypta_id_table=config.VERTICALS_NO_MULTI_PROFILE,  
crypta_id_to_hhid_table=config.HOUSEHOLD_CRYPTA_ID_TO_HHID,  
yandexuid_to_hhid_table=config.HOUSEHOLD_REVERSED_TABLE,  
hhid_to_yandexuid_table=config.HOUSEHOLD_ENRICH_TABLE,  
output_table=output_path,
```

Profiles integrate biometric data, most likely from smart speakers that use Yandex's Alice smart assistant



The image shows a screenshot of a Yandex website page. At the top, there is a red navigation bar with the Yandex logo on the left and menu items: COMPANY, JOBS, FOR DEVELOPERS, FOR ADVERTISERS, and FO. Below this is a white navigation bar with links: About, History, Privacy, Press Releases (highlighted in red), Blog, and Contact. The main content area has a breadcrumb trail: Press Releases / 2022 / followed by the article title: Yandex Launches Smart Devices With Alice in Uzbekistan. The article text begins with 'Internet, November 22, 2022. Uzbekistan's local Stations with Alice. Upon purchasing a smart sp' and is followed by a large image of three smart speakers (blue, black, and red) on a white surface. Below the image, the text continues: 'Stations are smart speakers with Alice. A single Yandex Plus subscription allows you to play music, podcasts, or playlists with personal recommendations on Yandex Station. Alice, the voice assistant on board, can entertain children with an educational game or compose a fairy tale together. She will tell you about the weather or remind you to buy groceries. Alice loves talking and will hold a conversation with ease: currently, in Russian only.'

Yandex

COMPANY JOBS FOR DEVELOPERS FOR ADVERTISERS FO

About History Privacy **Press Releases** Blog Contact

Press Releases / 2022 /

Yandex Launches Smart Devices With Alice in Uzbekistan

Internet, November 22, 2022. Uzbekistan's local Stations with Alice. Upon purchasing a smart sp

Stations are smart speakers with Alice. A single Yandex Plus subscription allows you to play music, podcasts, or playlists with personal recommendations on Yandex Station. Alice, the voice assistant on board, can entertain children with an educational game or compose a fairy tale together. She will tell you about the weather or remind you to buy groceries. Alice loves talking and will hold a conversation with ease: currently, in Russian only.

Possible Children by Voice

```
profile > runners > segments > lib > coded_segments > children_age_segment_clarification.py > ...
12
13   clarify_children_yql_template = """
14   $possible_children_by_voice = (
15     SELECT `uuid`, TableName() AS `date`, '0_12' AS segment_name
16     FROM RANGE(`{biometry_folder}`, `{biometry_first_date}`, `{biometry_last_date}`)
17     WHERE bio_child > 0.8
18   );
19
20   $possible_children_by_voice = (
21     SELECT DISTINCT `uuid`, `date`, segment_name
22     FROM $possible_children_by_voice
23   );
24
25   $possible_children_by_voice = (
26     SELECT `uuid`, segment_name
27     FROM $possible_children_by_voice
28     GROUP BY `uuid`, segment_name
29     HAVING COUNT(*) >= 2
30   );
31
32   $sources_new_age = (
33     SELECT matching.cryptaId AS cryptaId,
34     CASE
35       WHEN socdem_storage.birth_date > '{thirteenth_birthday}' THEN '0_12'
36       WHEN '{thirteenth_birthday}' >= socdem_storage.birth_date AND
37         socdem_storage.birth_date > '{eighteenth_birthday}' THEN '13_17'
38       ELSE '18_99'
39     END AS segment_name
40     FROM `{socdem_storage_table}` AS socdem_storage
41     INNER JOIN `{id_to_crypta_id_table}` AS matching
42     ON socdem_storage.id == matching.id AND socdem_storage.id_type == matching.id_type
43     WHERE socdem_storage.birth_date is not Null
44     UNION ALL
45     SELECT matching.cryptaId AS cryptaId, biometry.segment_name AS segment_name
46     FROM $possible_children_by_voice AS biometry
47     INNER JOIN `{id_to_crypta_id_table}` AS matching
48     ON biometry.`uuid` == matching.id
49     WHERE matching.id_type == 'uuid'
```

UI for Infographics Card

```
27
28 const marriedText = convertMarriedToSingleText(exactDemographics.gender, married);
29 const incomeText = convertIncomeSegmentToText(exactDemographics.income);
30 const hasChildrenText = convertHasChildrenToText(hasChildren);
31
32 return (
33   <div className="BasicInfoGraphics">
34     <img alt="" className="BasicInfoGraphics-Image" src={images[exactDemographics.gender]} />
35     <div className="BasicInfoGraphics-Bubble BasicInfoGraphics-Bubble_family">{marriedText}</div>
36     <div className="BasicInfoGraphics-Bubble BasicInfoGraphics-Bubble_income">{incomeText}</div>
37     <div className="BasicInfoGraphics-Bubble BasicInfoGraphics-Bubble_children">{hasChildrenText}</div>
38     <div className="BasicInfoGraphics-Interest BasicInfoGraphics-Interest_first">
39       <div className="BasicInfoGraphics-InterestIcon"
40         style={{ backgroundImage: `url(${interestIcons[0]})` }} />
41     </div>
42     <div className="BasicInfoGraphics-Interest BasicInfoGraphics-Interest_second">
43       <div className="BasicInfoGraphics-InterestIcon"
44         style={{ backgroundImage: `url(${interestIcons[1]})` }} />
45     </div>
46     <div className="BasicInfoGraphics-Interest BasicInfoGraphics-Interest_third">
47       <div className="BasicInfoGraphics-InterestIcon"
48         style={{ backgroundImage: `url(${interestIcons[2]})` }} />
49     </div>
50   </div>
51 );
52 }
```

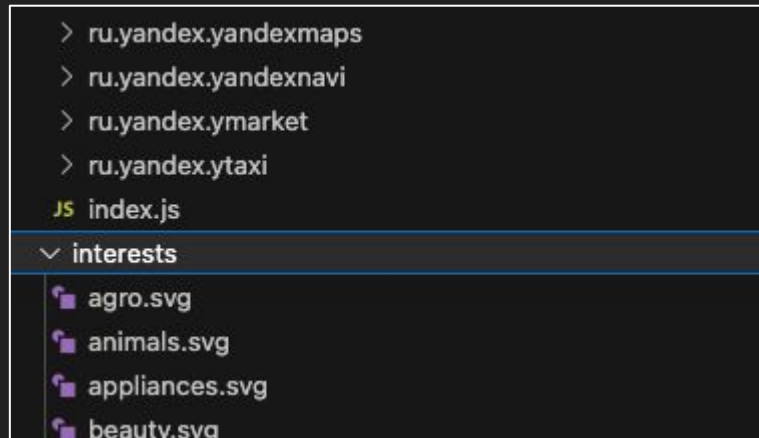
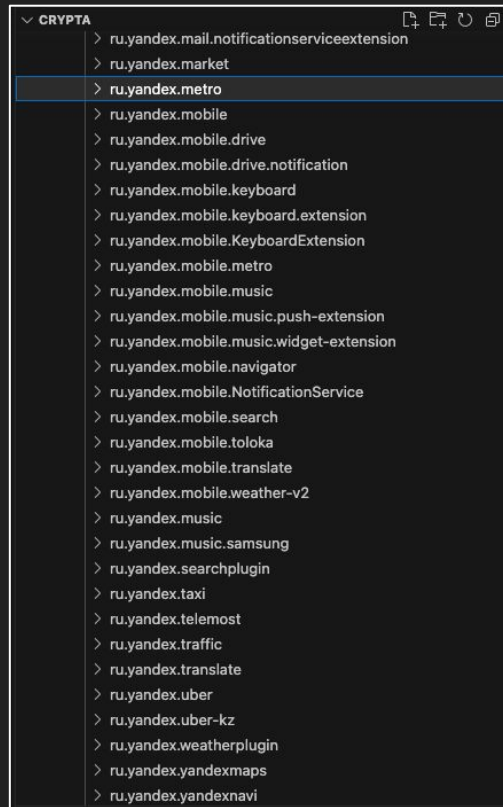
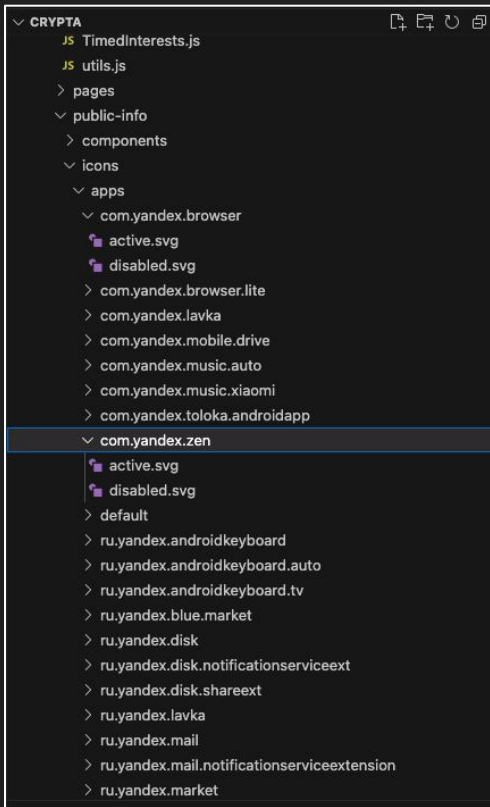
- icons
 - apps
 - interests
 - agro.svg
 - animals.svg
 - appliances.svg
 - beauty.svg
 - business.svg
 - clothes.svg
 - construction.svg
 - education.svg
 - electronics.svg
 - entertainments.svg
 - family.svg
 - finance.svg
 - food.svg
 - gifts.svg
 - JS index.js
 - job.svg
 - realty.svg
 - rest.svg
 - sport.svg
 - stationery.svg
 - telecom.svg
 - transport.svg

Search Profile by ID

```
web > portal > src > graph > search > JS SearchPanel.js > ...
112     return (
113         <div key={"inputs-" + suffix} className="input-group">
114             <div style={{display: showIdInput ? "block" : "none"}}>
115                 <ValueInput
116                     key="id_value"
117                     placeholder="[id value]"
118                     value={parameters["idValue" + suffix]}
119                     onChange={changeParameter("uid" + suffix)}
```

```
return (
  experiments.status !== 403 && (
    <div className="experiments-bar">
      <div className="experiments-select-uid-type">
        <RadioButton
          value={activeUid}
          size="s"
          view="default"
          className="select-sorting"
          onChange={(event) => selectUidType(event.target.value)}
          options={[
            { value: "uid", children: t("by") + " yandexuid" },
            { value: "cryptaId", children: t("by") + " CryptaID" },
          ]}
        />
      </div>
    </div>
  )
)
```


UI - Available App Icons



Ids Associated with Social Media Accounts

```
web > portal > src > public-info > sections > GraphSection > JS GraphSection.js > GraphSection >
1 import React, { useEffect, useMemo, useState } from "react";
2 import { useSelector } from "react-redux";
3 import { getPublicGraph, getPublicGraphLoading } from "../../store/selectors";
4 import { Graph, GraphSkeleton } from "../../components/Graph/Graph";
5 import { Section } from "../../components/Section/Section";
6 import { getServiceIcon } from "../../icons/services";
7 import { getAppIcon } from "../../icons/apps";
8
9 import "./GraphSection.scss";
10
11 import noData from "./no-data.svg";
12
13 const IMAGE_SIZE_XS = 12;
14 const IMAGE_SIZE_S = 36;
15 const IMAGE_SIZE_M = 56;
16 const IMAGE_SIZE_L = 80;
17
18 const NODE_MAPPING = {
19   email: {
20     imageSize: IMAGE_SIZE_M,
21     imageHref: "mail",
22   },
23   yandexuid: {
24     imageSize: IMAGE_SIZE_XS,
25     imageHref: "yandexuid",
26   },
27   idfa: {
28     imageSize: IMAGE_SIZE_L,
29     imageHref: "ios",
30   },
31   gaid: {
32     imageSize: IMAGE_SIZE_L,
33     imageHref: "android",
34   },
35   oaid: {
36     imageSize: IMAGE_SIZE_L,
37     imageHref: "android",
38   },
39   login: {
40     imageSize: IMAGE_SIZE_M,
41     imageHref: "key",
42   },
43   puid: {
44     imageSize: IMAGE_SIZE_M,
45     imageHref: "key",
46   },
47   instagram_login: {
48     imageSize: IMAGE_SIZE_M,
49     imageHref: "instagram",
50   },
51   instagram_id: {
52     imageSize: IMAGE_SIZE_M,
```

```
web > portal > src > public-info > sections > GraphSection > JS GraphSection.js > GraphSection > useEffec
51   },
52   instagram_id: {
53     imageSize: IMAGE_SIZE_M,
54     imageHref: "instagram",
55   },
56   fb_id: {
57     imageSize: IMAGE_SIZE_M,
58     imageHref: "facebook",
59   },
60   ok_id: {
61     imageSize: IMAGE_SIZE_M,
62     imageHref: "ok",
63   },
64   vk_id: {
65     imageSize: IMAGE_SIZE_M,
66     imageHref: "vk",
67   },
68   vk_name: {
69     imageSize: IMAGE_SIZE_M,
70     imageHref: "vk",
71   },
72   kp_id: {
73     imageSize: IMAGE_SIZE_M,
74     imageHref: "kinopoisk",
75   },
76 }
77
78 function getNodeMapping(item) {
79   if (item.idType === 'uuid') {
80     // App
81     return { imageHref: item.icon, imageSize: IMAGE_SIZE_S };
82   }
83
84   return NODE_MAPPING[item.icon] ?? { imageHref: "default", imageSize: IMAGE_SIZE_XS };
85 }
86
87 function getImage(item) {
88   const disabled = !item.isActive;
89
90   if (item.idType === 'uuid') {
91     return getAppIcon(item.imageHref, disabled)
92       .catch(() => getAppIcon("default", disabled));
93   }
94
95   return getServiceIcon(item.imageHref, disabled)
96     .catch(() => getServiceIcon("default", disabled));
97 }
```


Matcher

Matcher

- ∨ matcher
 - > bin
 - > bundle
 - ∨ lib
 - > config
 - ∨ matchers
 - > base_matcher
 - > beeline_matcher
 - > er_telecom_matcher
 - > intentai_matcher
 - > mts_matcher
 - > rostelecom_matcher
 - ≡ ya.make
 - ⊕ parser.cpp
 - ⊕ parser.h

Rostelecom Matcher

```
ext_fp > matcher > lib > matchers > rostelecom_matcher > C rostelecom_matcher.cpp
22 TConnection TRostelecomMatcher::MakeConnection(const TFpEvent& event) {
23     return {
24         .Ip = event.GetIp(),
25         .Port = event.GetPort(),
26         .Timestamp = event.GetUnixtime(),
27         .Domain = NMcDomain::GetMcDomainForRostelecom(event.GetDuid()),
28     };
29 }
30
31 void TRostelecomMatcher::AddConnection(const TFpEvent& event) {
32     auto connection = MakeConnection(event);
33
34     Stats.Count->Add("events.incoming.rostelecom.count");
35     Request += TStringBuilder() << connection.Ip << '\t'
36         << connection.Port << '\t'
37         << connection.Timestamp << '\t'
38         << connection.Domain << '\n';
39 }
40
41 TMatches TRostelecomMatcher::GetMatches() {
42     if (Request.length() == 0) {
43         return TMatches();
44     }
45     const auto& requestId = CreateGuidAsString();
46     Log->info("Rostelecom request {} body:\n{}", requestId, Request);
47
48     NNeh::TMessage message(GetApiUrl(), "");
49     Y_ENSURE(NNeh::NHttp::MakeFullRequest(message, "", Request, "text/plain"), "Failed to build request to Rostelecom API");
50
51     Stats.Count->Add("api.calls.rostelecom.count");
52     const auto& resp = MakeRequest(Client, message, TDuration::Milliseconds(Config.GetApiCallTimeoutMs()), "Rostelecom", requestId, Log);
53
54     return ParseResponse(resp->Data);
55 }
56
57 TString TRostelecomMatcher::GetApiUrl() const {
58     return "post://" + Config.GetApiUrl();

```

Rostelecom Matcher

```
ext_fp > matcher > lib > matchers > rostelecom_matcher > rostelecom_matcher.cpp
22   TConnection TRostelecomMatcher::MakeConnection(const TFpEvent& event) {
23       return {
24           .Ip = event.GetIp(),
25           .Port = event.GetPort(),
26           .Timestamp = event.GetUnixtime(),
27           .Domain = NMcDomain::GetMcDomainForRostelecom(event.GetDuid()),
28       };
29   }
```

Rostelecom Matcher

```
41 TMatches TRostelecomMatcher::GetMatches() {  
42     if (Request.length() == 0) {  
43         return TMatches();  
44     }  
45     const auto& requestId = CreateGuidAsString();  
46     Log->info("Rostelecom request {} body:\n{}", requestId, Request);  
47  
48     NNeh::TMessage message(GetApiUrl(), "");  
49     Y_ENSURE(NNeh::NHttp::MakeFullRequest(message, "", Request, "text/plain"), "Failed to build request to Rostelecom API");  
50  
51     Stats.Count->Add("api.calls.rostelecom.count");  
52     const auto& resp = MakeRequest(Client, message, TDuration::Milliseconds(Config.GetApiCallTimeoutMs()), "Rostelecom", requestId, Log);  
53  
54     return ParseResponse(resp->Data);  
55 }
```

Test Result Data

```
ext_fp > matcher > bin > test > canondata > {} result.json > [ ] test_matcher.test_matcher > { 1
  2 {
  3   "test_matcher.test_matcher": [
  4     {
  5       "duid": "16999999761000006",
  6       "ext_id": "fake_ertelecom_id_for_5.3.100.0",
  7       "ext_source": "ertelecom",
  8       "hitlogid": "100506",
  9       "ip": "5.3.100.0",
10       "log_type": "bs-watch-log",
11       "logid": 0,
12       "original_domain": "domain-6.ru",
13       "port": "5555",
14       "rtmr_timestamp": "1699999977",
15       "unixtime": "1699999970",
16       "user_agent": "Mozilla/5.0 (Windows NT PYQ)",
17       "watchid": "20000000000000006",
18       "yuid": "10061699999976"
19     },
20     {
21       "duid": "16999999861000016",
22       "ext_id": "mts_id_for_160.1.2.4",
23       "ext_source": "mts",
24       "hitlogid": "100516",
25       "ip": "160.1.2.4",
26       "log_type": "bs-watch-log",
27       "logid": 0,
28       "original_domain": "domain-16.ru",
29       "port": "4444",
30       "rtmr_timestamp": "1699999987",
31       "unixtime": "1699999970",
32       "user_agent": "Mozilla/5.0 (Windows NT PYQ)",
33       "watchid": "20000000000000016",
34       "yuid": "10161699999986"
35     },
36     {
37       "duid": "16999999931000003",
38       "ext_id": "fake_ertelecom_id_for_5.3.62.0",
39       "ext_source": "ertelecom",
40       "hitlogid": "100503",
41       "ip": "5.3.62.0",
42       "log_type": "bs-watch-log",
43       "logid": 0,
44       "original_domain": "domain-3.ru",
45       "port": "2222",
46       "rtmr_timestamp": "1699999994",
47       "unixtime": "1699999990",
48       "user_agent": "Mozilla/5.0 (Windows NT PYQ)",
49       "watchid": "20000000000000003",
50       "yuid": "10031699999993"
51     }
52   ]
53 }
```

Roadmap

- Background on Yandex Leak
- Dive into code:
 - What data Yandex is collecting
 - What Yandex is doing with that data
 - Who Yandex is sharing that data with
- Conclusions and wrap up
- Q&A

Conclusion

Wrap Up

- Yandex has access to a broad international reach of data and it has been evasive about what it can do with that data
- A small amount of data can say a lot when it is matched to entries from a company's other data sources and analyzed
- Yandex has code to sync some of its data with a Russian-state owned entity

Takeaways

- Anonymization is very easily undone when data gets combined with pools from other sources that may contain identifying data
- Pay attention to who runs your SDKs, what data points they collect, and where they send your user data.
- Who gets access to a company's user data when its assets are sold, the geopolitical climate changes, or a government tightens its control?

Q&A

Link to Write Up:
<https://bit.ly/455utBP>