

What the Yandex Leak Tells Us About How Big Tech Uses Your Data

Kaileigh McCrea, Privacy Engineer @ Confiant

August 2023



What we're talking about

Yandex

The screenshot shows a post on the BreachForums forum. The title is "yandex git sources" and it was posted by "breachforums" on Wednesday January 25, 2023 at 03:48 PM. The post content is a link to a GitHub repository containing source code for Yandex services. The user has 1 post, 1 thread, joined in Jan 2023, and a reputation of 0. There are options to reply, quote, or report the post.

YANDEX SERVICES SOURCE CODE LEAK

SHORT OVERVIEW OF BREACH CONTENTS

PUBLISHED THU, JAN 26, 2023 BY ARSENIY SHESTAKOV

The screenshot shows a news article from BleepingComputer. The headline is "Yandex denies hack, blames source code leak on former employee". The article is by Bill Toulas and was published on January 26, 2023, at 09:44 AM. The BleepingComputer navigation bar includes links for NEWS, DOWNLOADS, VPNs, VIRUS REMOVAL GUIDES, TUTORIALS, and DEALS. There is also a search bar and social media links for Facebook, Twitter, and YouTube.

Roadmap

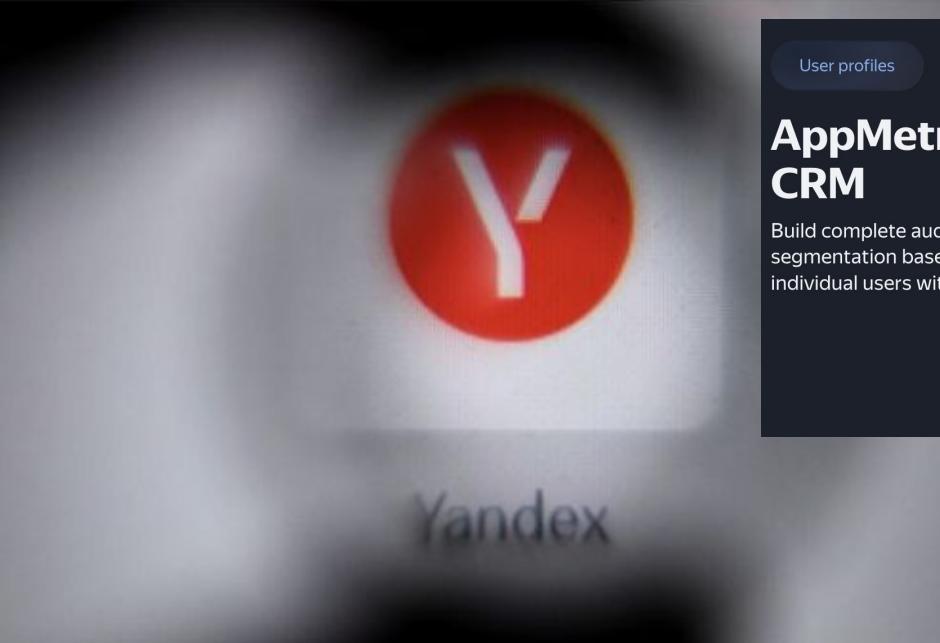
- Background on Yandex Leak
- Dive into code:
 - What data Yandex is collecting
 - What Yandex is doing with that data
 - Who Yandex is sharing that data with
- Conclusions and wrap up
- Q&A

Yandex: A Drama

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



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.

♂ 18-24
📍 ⚽

Today

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

2022

September

November

December

January

Russia invades Ukraine

Alarm About Sending Data to Russian Servers

February

March

April

““Yandex has acknowledged its software collects “device, network and IP address” information ... but it called this data “non-personalised and very limited”. It added: “Although theoretically possible, in practice it is extremely hard to identify users based solely on such information collected. Yandex definitely cannot do this.””

– *Financial Times*

Featured Article

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

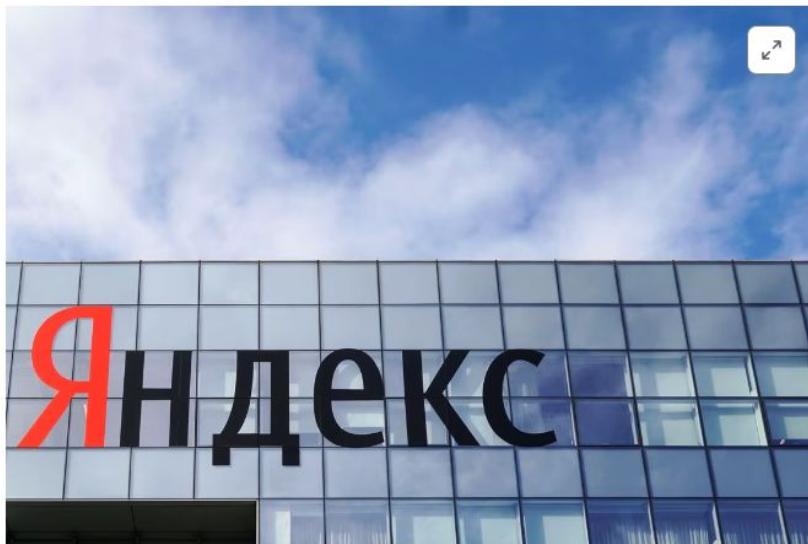
 Comment



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



≡ NEWS UKRAINE WAR BUSINESS OPINION ARTS AND LIFE PODCASTS NEWSLETTERS ARCHIVE

Russian Billionaires Line Up to Buy Yandex – Reports

May 4, 2023



MOST READ JUST IN

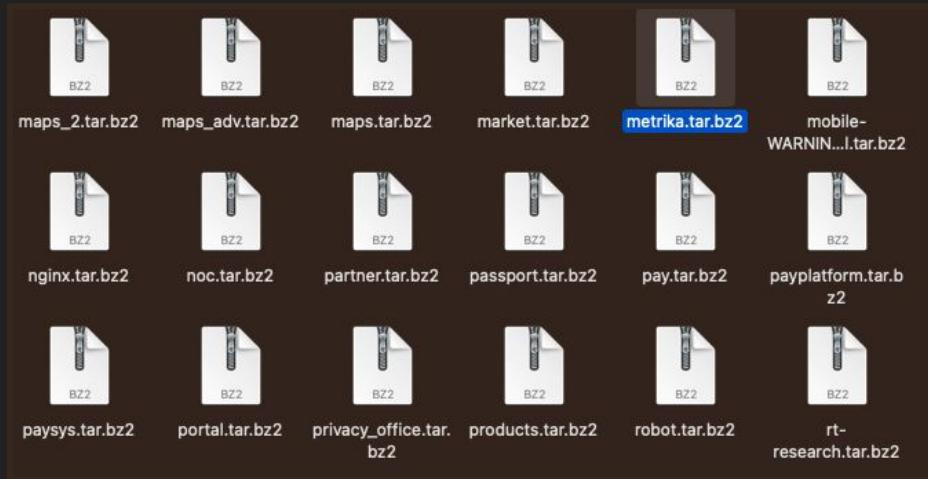
- 1 [OVERNIGHT STRIKE](#)
Russia Says Ukrainian Droned Moscow, Crimea
- 2 [NO PASSAGE](#)
Russia Blocks Cargo Ship On Traces'
- 3 [POLITICAL PRISONER](#)
Navalny Ally Jailed 9 Years
- 4 [MORE MANPOWER](#)
Russia Raises Upper-Age Limit Reservists
- 5 [MONEY DRAIN](#)
Russia Loses Record \$253B Capital Flight

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

[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 Metrica

[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.yql

```
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_PhysicalCellsIDs`,  
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` ,
```

analytics > appmetrica-location-log-anonymizer > convert_log.yql

```
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     `timestamp` ,
```

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`  
145 `ReceiveTimestamp`  
146 `RequestID`
```

Wifi Fields Collected by AppMetrica

```
147     `SendTimestamp` ,  
148     `Wifi_AreConnected` ,  
149     `Wifi_LastVisibleTimeOffset` ,  
150     `Wifi_Macs` ,  
151     `Wifi_SignalsStrengths` ,  
152     `Wifi_Ssids` ,  
153     `loafeller index bucket` .
```

Those Same 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_ssids_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         }
```

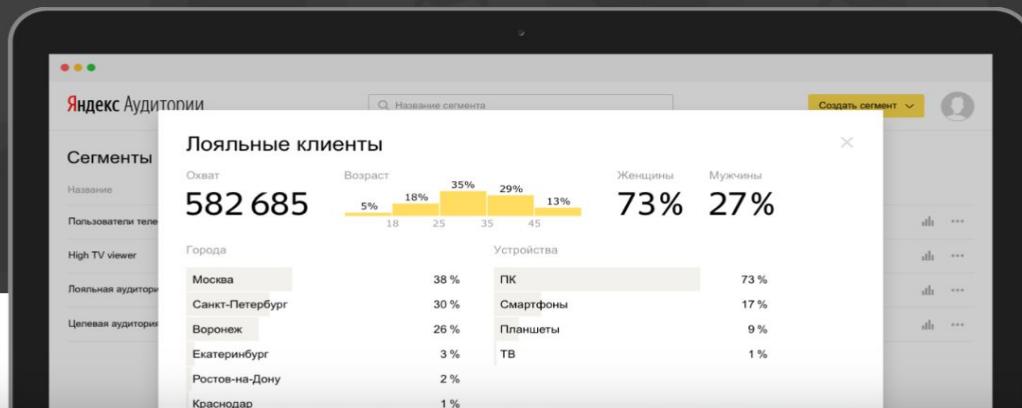
Socio-Demographic Attributes for DevID Being Updated

```
core > programs > socdem-updated-mobile > src > UserIdAndInfoParser.cpp > ...
45     t"0_1", AgeIntervalsCrypt::LessThan18},
46     {"18_24", AgeIntervalsCrypt::Between18and24},
47     {"25_34", AgeIntervalsCrypt::Between25and34},
48     {"35_44", AgeIntervalsCrypt::Between35and44},
49     {"45_54", AgeIntervalsCrypt::Between45and54},
50     {"55_99", AgeIntervalsCrypt::MoreThan55}
51 };
52 ::setValue(value, exact_socdem_node, key, json_keys_to_ages_intervals);
53 }
54
55 void UserIdAndInfoParser::setValue(
56     SexTypesCrypt & value,
57     const NYT::TNode & exact_socdem_node,
58     const TString key)
59 {
60     static const std::map<TString, SexTypesCrypt> json_keys_to_sex_types =
61     {
62         {"f", SexTypesCrypt::Female},
63         {"m", SexTypesCrypt::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



Crypta

The screenshot shows a dark-themed web page for Yandex. At the top, there is a red header bar with the Yandex logo on the left and navigation links for COMPANY, JOBS, FOR DEVELOPERS, FOR ADVERTISERS, and FOR INVESTORS. To the right of these links is a vertical three-line menu icon. Below the red bar is a light gray navigation bar containing links for About, History, Privacy, Press Releases, Blog, Contact, and a magnifying glass icon for search. The main content area has a white background. At the top of this area, there is a breadcrumb trail with 'Technologies /' followed by a large, bold, black title 'Crypta'. Below the title, there is a paragraph of text describing the technology.

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	INTERESTS	INTERESTS
> test		
age_segment_18_20.py	kinopoisk_movie_watchers.py	mobile_operators_users_by_prefix.py
alice_users.py	kz_users.py	mobile_operators_users.py
apartment_room_number.py	laptop_users.py	multidevice_puid.py
apps_users.py	logged_in_for_plus.py	multidevice.py
artists.py	longterm_interest_mobile_gamers.py	music_genres_listeners.py
auto_interactions.py	loyal_to_launcher_install.py	nestle_regions.py
avia_travellers.py	macos_users.py	phone_buyers.py
bank_cards.py	mail_data.py	phone_owners.py
bought_two_tickets.py	manufacturer_phone_owners.py	phone_with_esim_owners.py
business_travellers.py	marketplaces_ltv_users.py	potential_aon_android_users.py
children_age_segment_clarification.py	mobile_gamers.py	potential_aon_ios_users.py
compulsory_auto_insurance.py	mobile_operators_users_by_prefix.py	preinstalled_apps.py
connection_type.py	mobile_operators_users.py	prism.py
console_gamers.py	multidevice_puid.py	proleads.py
contest.py	multidevice.py	realty_interactions.py
devices_without_google_services.py	music_genres_listeners.py	recent_passport_accounts.py
digital_viewers.py	nestle_regions.py	score_users_for_telephony.py
direct_clients_by_industry.py	phone_buyers.py	searched_for_phone_numbers.py
direct_product_users.py	phone_owners.py	searched_radisson_on_maps.py
disk_users.py	phone_with_esim_owners.py	seo_specialists.py
ecommerce_owners.py	potential_aon_android_users.py	seo_users.py
edadeal_offline_purchases_lal.py	potential_aon_ios_users.py	smart_gadgets_customers.py
expensive_car_customers.py	preinstalled_apps.py	smokers.py
film_lovers_by_genres.py	prism.py	summer_residents.py
gas_stations.py	proleads.py	travellers.py
industryRepresentatives.py	realty_interactions.py	video_bloggers.py
kbt_customers.py	recent_passport_accounts.py	want_to_change_the_provider.py
kfc_visitors.py	score_users_for_telephony.py	webmaster.py
kinopoisk_logins.py	searched_for_phone_numbers.py	widgets.py
kinopoisk_movie_watchers.py	searched_radisson_on_maps.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', 'crypto_id'),
103                     matching_gaid=get_matching_table('gaid', 'crypto_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             )
```

Basic Example of Household Details

```
graph > metrics > household > query.sql
42  DEFINE SUBQUERY $hh_size_by_yuid($title) AS
43    SELECT ($title || key) AS key, COUNT(1) AS hh_size_by_yuids
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 != "";
    }
```

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             ],
96         
```

Appmetrica Data Being Used to Identify 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 Identify 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",
74                 "computerid",
75                 "cid",
76                 "suserid",
77             ]
78         ).SerializeToString()
79
```



Extracting
multiple types
of identifiers

Passport

The screenshot shows a web browser window with the URL <https://passport.yandex.com/registration?retpath=https://audience.yandex.com/>. The browser's address bar and various tabs are visible at the top. On the left side of the page, there is a large red Yandex logo and several smaller icons representing different services: a location pin, a red Y logo, a document with an 'A' and Chinese characters, an envelope, and a blue globe.

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
```

Geo Graphs

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

```
graph > fuzzy > lib > tasks > sources > geo > C geo_operations.h
61     for (; input->IsValid(); input->Next()) {
62         const auto& row = input->GetRow();
63         if (not IsRowValid(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         }
    }
```

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 to Generate Segments about Households with Children

```
profile > runners > segments > lib > coded_segments > with_children_by_ages.py > ...
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
def prepare_with_children_by_app(self, with_children_by_app_table):
    app_segment_name_tuples = []

    for segment_name, apps in SEGMENT_APP_DICT.iteritems():
        for app in apps:
            app_segment_name_tuples.append('AsTuple("{}", "{}").format(app, segment_name)')

    self.yql.query(
        app_metric_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_metric_table:
        self.yt.run_map(
            extract_children_from_taxi,
            inputs['TaxiData'].table,
            taxi_puid_table,
        )

        self.prepare_with_children_by_app(app_metric_table)

        self.yql.query(
            with_children_query_template.format(
                metrics_table=inputs['ProcessedMetrics'].table,
                reqans_table=inputs['ProcessedReqans'].table,
                app_metric_table=app_metric_table,
                taxi_data_table=taxi_puid_table,
                id_to_crypt_id_table=config.VERTICES_NO_MULTI_PROFILE,
                crypt_id_to_hhid_table=config.HOUSEHOLD_CRYPT_ID_TO_HHID,
                yandexuid_to_hhid_table=config.HOUSEHOLD_REVERSED_TABLE,
                hhid_to_yandexuid_table=config.HOUSEHOLD_ENRICH_TABLE,
                output_table=output_path,
            ),
            transaction=self.transaction,
        )
```

Appmetrica and Taxi Data Being Used to Generate Segments about Households with Children

```
456
457     def build_segment(self, inputs, output_path):
458         with self.yt.TempTable() as taxi_puid_table, \
459             self.yt.TempTable() as app_metrica_table:
460             self.yt.run_map(
461                 extract_children_from_taxi,
462                 inputs['TaxiData'].table,
463                 taxi_puid_table,
464             )
465
466             self.prepare_with_children_by_app(app_metrica_table)
467
```

ID Mapping Associations

```
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.VERTICES_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,  
    ),  
    transaction=self.transaction,  
)
```

Profiles Integrate Biometric Data, Most Likely from Smart Speakers That Use Yandex's Alice Smart Assistant

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 retailers Stations with Alice. Upon purchasing a smart speaker, us

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.cryptoId AS cryptoId,
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_crypto_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.cryptoId AS cryptoId, biometry.segment_name AS segment_name
46         FROM $possible_children_by_voice AS biometry
47         INNER JOIN `{id_to_crypto_id_table}` AS matching
48         ON biometry.`uuid` == matching.id
49         WHERE matching.id_type == 'uuid'
50     );

```

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

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 > useEffect
  ...
 51   instagram_id: {
 52     imageSize: IMAGE_SIZE_M,
 53     imageHref: "instagram"
 54   },
 55   fb_id: {
 56     imageSize: IMAGE_SIZE_M,
 57     imageHref: "facebook"
 58   },
 59   ok_id: {
 60     imageSize: IMAGE_SIZE_M,
 61     imageHref: "ok"
 62   },
 63   vk_id: {
 64     imageSize: IMAGE_SIZE_M,
 65     imageHref: "vk"
 66   },
 67   vk_name: {
 68     imageSize: IMAGE_SIZE_M,
 69     imageHref: "vk"
 70   },
 71   kp_id: {
 72     imageSize: IMAGE_SIZE_M,
 73     imageHref: "kinopoisk"
 74   }
 75 }

function getNodeMapping(item) {
  if (item.idType === 'uuid') {
    // App
    return { imageHref: item.icon, imageSize: IMAGE_SIZE_S };
  }

  return NODE_MAPPING[item.icon] ?? { imageHref: "default", imageSize: IMAGE_SIZE_XS };
}

function getImage(item) {
  const disabled = !item.isActive;

  if (item.idType === 'uuid') {
    return getAppIcon(item.imageHref, disabled)
      .catch(() => getAppIcon("default", disabled));
  }

  return getServiceIcon(item.imageHref, disabled)
    .catch(() => getServiceIcon("default", disabled));
```

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
      C parser.h
```

Rostelecom Matcher

```
ext_fp > matcher > lib > matchers > rostelecom_matcher > rostelecom_matcher.cpp
 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();
 59 }
 60
 61 TMatches TRostelecomMatcher::ParseResponse(const TString& response) {
 62     TStringInput stringInput(response);
 63
 64     TString line;
 65     TConnection connection;
 66     TString extId;
 67     TMatches matches;
 68
 69     while (stringInput.ReadLine(line)) {
 70         Split(TStringBuf(line), '\t', connection.Ip, connection.Port, connection.Timestamp, connection.Domain, extId);
 71         matches[connection] = TMatchResult{.Status = TMatchResult::EStatus::Found, .ExtId = extId};
 72     }
 73
 74     return matches;
 75 }
```

Rostelecom Matcher

```
ext_fp > matcher > lib > matchers > rostecom_matcher > rostecom_matcher.cpp
31 void TRostelecomMatcher::AddConnection(const TFpEvent& event) {
32     auto connection = MakeConnection(event);
33
34     Stats.Count->Add("events.incoming.rostecom.count");
35     Request += TStringBuilder() << connection.Ip << '\t'
36                                     << connection.Port << '\t'
37                                     << connection.Timestamp << '\t'
38                                     << connection.Domain << '\n';
39 }
40
```

Rostelecom Matcher

```
61 TMatches TRostelecomMatcher::ParseResponse(const TString& response) {
62     TStringInput stringInput(response);
63
64     TString line;
65     TConnection connection;
66     TString extId;
67     TMatches matches;
68
69     while (stringInput.ReadLine(line)) {
70         Split(TStringBuf(line), '\t', connection.Ip, connection.Port, connection.Timestamp, connection.Domain, extId);
71         matches[connection] = TMatchResult{.Status = TMatchResult::EStatus::Found, .ExtId = extId};
72     }
73
74     return matches;
75 }
```

Test Result Data

ext_fp > matcher > bin > test > canondata > {} result.json > []test_matcher.test_matcher > {} 1

```
1  {
2      "test_matcher.test_matcher": [
3          {
4              "duid": 16999999761000006,
5              "ext_id": "fake_ertelecom_id_for_5.3.100.0",
6              "ext_source": "ertelecom",
7              "hitlogid": 100506,
8              "ip": "5.3.100.0",
9              "log_type": "bs-watch-log",
10             "logid": 0,
11             "original_domain": "domain-6.ru",
12             "port": 5555,
13             "rtmr_timestamp": 1699999977,
14             "unixtime": 1699999970,
15             "user_agent": "Mozilla/5.0 (Windows NT [PYC])",
16             "watchid": 200000000000000006,
17             "yuid": 1006169999976
18         },
19         {
20             "duid": 16999999861000016,
21             "ext_id": "mts_id_for_160.1.2.4",
22             "ext_source": "mts",
23             "hitlogid": 100516,
24             "ip": "160.1.2.4",
25             "log_type": "bs-watch-log",
26             "logid": 0,
27             "original_domain": "domain-16.ru",
28             "port": 4444,
29             "rtmr_timestamp": 1699999987,
30             "unixtime": 1699999970,
31             "user_agent": "Mozilla/5.0 (Windows NT [PYC])",
32             "watchid": 200000000000000016,
33             "yuid": 1016169999986
34         },
35         {
36             "duid": 16999999931000003,
37             "ext_id": "fake_ertelecom_id_for_5.3.62.0",
38             "ext_source": "ertelecom",
39             "hitlogid": 100503,
40             "ip": "5.3.62.0",
41             "log_type": "bs-watch-log",
42             "logid": 0,
43             "original_domain": "domain-3.ru",
44             "port": 2222,
45             "rtmr_timestamp": 1699999994,
46             "unixtime": 1699999990,
47             "user_agent": "Mozilla/5.0 (Windows NT [PYC])",
48             "watchid": 200000000000000003,
49             "yuid": 1003169999993
50         }
51     ]
52 }
```

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

What they have:



WiFi SSID



Device information



Location



Search history

What they can determine:



Household composition



Age



Income



Gender



Biometric voice data



E-Mail



Phone number



Children's ages



Interests



Home location



Birth date



App usage data



Service usage data



Physical neighbors



Travel plans/
history



Switching providers



Similar users



Proximity to
local businesses



Smoker



Music
preferences



Home location



Work location



Movie/music
history



E-Mail data

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?



Link to Write Up:
(available August 10th)
bit.ly/455utBP

Q & A



Link to Slides:
(available August 10th)
bit.ly/3qail3p