

**U.S. Department of Commerce
U.S. Patent and Trademark Office**



**Privacy Threshold Analysis
for the
Open Data-Big Data Master System (OD-BD MS)**

U.S. Department of Commerce Privacy Threshold Analysis

USPTO Open Data-Big Data Master System (OD-BD MS)

Unique Project Identifier: PTOC-034-00

Introduction: This Privacy Threshold Analysis (PTA) is a questionnaire to assist with determining if a Privacy Impact Assessment (PIA) is necessary for this IT system. This PTA is primarily based from the Office of Management and Budget (OMB) privacy guidance and the Department of Commerce (DOC) IT security/privacy policy. If questions arise or further guidance is needed in order to complete this PTA, please contact your Bureau Chief Privacy Officer (BCPO).

Description of the information system and its purpose: *Provide a general description (in a way that a non-technical person can understand) of the information system that addresses the following elements:*

The E-Government Act of 2002 defines “information system” by reference to the definition section of Title 44 of the United States Code. The following is a summary of the definition: “Information system” means a discrete set of information resources organized for the collection, processing, maintenance, use, sharing, dissemination, or disposition of information. See: 44 U.S.C. § 3502(8).

The Open Data/Big Data (OD-BD) master system consists of subsystems that support the Big Data Portfolio. OD-BD MS resides on the UACS platform, which employs Infrastructure as a Service (IaaS) and Platform as a Service (PaaS) services from AWS and is located at USPTO Headquarters located at 600 Dulany Street, Alexandria, Virginia 22314 (“IT EAST” Environment). Subsystem details are provided below:

BDR: The Big Data Reservoir provides USPTO employees a Big Data platform in which they can view records and associated metadata in one location. The Big Data Reservoir (BDR) is a Hadoop Distributed File System (HDFS) infrastructure used to perform advanced analytics on disparate data sets consisting of structured and unstructured data in order to gain insights and develop models. System users are USPTO Internal Users. BDR is a large repository for structured and unstructured data. Models and algorithms are developed with the BDR data to provide insights to PTO executives. Dashboards, search functionality, and visualizations provide users the ability to view the BDR data.

- **BDR-TQR:** In addition to the BDR Portal, the BDR also provides the Trademark Quality Review (TQR) Portal. The TQR Portal provides quality reviewers with a centralized location to view the Dockets that are in the queue for review and additional features that include reviewing Trademark Review forms and completing necessary actions, final and non-final. System users are USPTO Internal Users.
- **BDR-CPC:** Cooperative Patent Classification (CPC) is used to automatically classify patent documents. Users can place input .csv file by using SFTP, which contains number of application IDs. By using this input file, BDR AI API gets the contractor data from the CPC OracleDB and machine data from the BDR AI for corresponding application IDs and stores the data in two csv files. Users can

compare contractor data and machine data, by giving application ID in the WEB. System users are USPTO Internal Users.

- **BDR-PTAB:** Patent Trial Appeal Board (PTAB) uses the BDR framework to gather data from PTAB E2E Oracle DB (PALMGP) and also from two OPSG REST APIs. Newly populated data in Oracle DB is collected by using Delta processing and stored in BDR HIVE/HDFS locations. The entire hive table's data is stored in SOLR index (Public), Elastic Index, and users can easily search the data based on a particular attribute. System users are USPTO Internal Users.

a) Whether it is a general support system, major application, or other type of system

There are multiple components to BDR; therefore, it would fall under an "Other" type of system:

- Advanced analytics infrastructure with a front-end user interface with dashboard, search, and visualization functionality.
- Trademark application to capture Quality Review information.

b) System location

USPTO Headquarters located at 600 Dulany Street Alexandria, Virginia 22314 ("IT EAST" Environment).

c) Whether it is a standalone system or interconnects with other systems (identifying and describing any other systems to which it interconnects)

The following table provides a list of applications that supply data to BDR:

| | |
|-----------------|--|
| BDSS | BDR gets Full text patents data, Full text grants data, Biblio Patents data, Biblio grants data, Patent Application PDFs and Grant Application PDFs from BDSS. After processing this data, Patents attributes will be pushed SOLR index <code>ibd_grantsv1</code> and Patents attributes will be pushed SOLR index <code>ibd_publicationsv1</code> . |
| TEAS | BDR pulls all Trademark application data from TEAS |
| PATI-CDC | BDR CPC interfaces with PATI-CDC to get meta data (documentCode, documentLocationURI, fileSize, fileSizeUnitCode, markupStandardCategory, documentLoadedDateTime, vendorName) for corresponding CPC Application numbers. |
| TRM | BDR-TQR interfaces with TRM to get PUBS and SOU data used to populate BDR-TQR review screen |
| CPC DB | BDR CPC interfaces with CPC DB to get contractor data (CPC codes). Users can compare this contractor data with machine data in CPC WEB by giving the corresponding application ID. |

| | |
|---------------------|--|
| OPSG | BDR interfaces with OPSG to get PTAB Appeal proceedings data to populate BDR-PTAB and DH-PTAB. BDR also retrieves Application IDs related to Office Actions for OA processing. |
| PTAB-E2E | BDR interfaces with PTAB E2E to get PTAB Trials, Appeals and Interferences data. After processing this data, internal data will be pushed to Elastic index (source BDR-PTAB), and public data will be pushed to SOLR indices and AWS S3 (Sources for DH-PTAB). |
| FAST2 | BDR-TQR interfaces with FAST2 to get tagged data of First Office Actions and Final Office Actions to identify form paragraphs used within an Office Action. This is used to populate BDR-TQR review screen. |
| PALM-EXPO | BDR uses application information and status codes pulled from PALM-GP tables to provide information on application attributes |
| PALM-INFRA | BDR uses application information and status codes pulled from PALM-GP tables to provide information on application attributes |
| PALM-PreExam | BDR uses application information and status codes pulled from PALM-GP tables to provide information on application attributes |
| RBAC | RBAC provides role based accessed control for BDR portal, infrastructure, and BDR services. |
| TMNG-CMS | BDR pulls mark image for display within BDR-TQR |
| P-ELP | BDR uses P-ELP CMS services to get Patent Office Action data – JSON metadata. In addition, BDR uses the P-ELP CMS services to retrieve patent Office Action XML files. |

d) The purpose that the system is designed to serve

The system is designed to serve as the enterprise platform for advanced analytics.

e) The way the system operates to achieve the purpose

BDR is a large repository for structured and unstructured data. Models and algorithms are developed with the BDR data to provide insights to PTO executives. Dashboards, search functionality, and visualizations provide users the ability to view the BDR data.

f) A general description of the type of information collected, maintained, used, or disseminated by the system

a. Patent Data

- i. Publicly disseminated data (PG Pubs/Grants, PTAB decisions)*
- ii. Patent Attributes (PALM)*
- iii. Patent Office Actions (P-ELP), PATI-CDC data*

- b. Trademark –
 - i. FAST – tagged paragraphs
 - ii. TRM – Trademark application information
 - iii. CMS – Trademark mark graphic
 - iv. BDR-TQR – TQR UI captures Trademark review information

g) *Identify individuals who have access to information on the system*

Users who have access to the system include Patent executives, Trademark reviewers, Data Scientists, and System Administrators/Operators.

h) *How information in the system is retrieved by the user*

BDR is a large repository for structured and unstructured data. There is the compute tier, where the data is loaded, compared for public versus private status, and analyzed according to data science principles. There is the analysis tier, where data scientists combine the real world problem solving techniques from Patent Examiners with the formulae and hypothesis of the Data Science field. The Visualization tier provides the users with a place to view the analysis and the underlying data that helps to create it. Finally, in the storage tier, the system retains raw, merged and transformed data, distinguishes between public and private Patent applications, and segregates them. Dashboards, search functionality, and visualizations provide users the ability to view the BDR data.

Developer Hub (DH) uses an N-tier architectural design pattern that separates the processing logic into distinct processing layers. The system is logically divided into six major subsystems:

1. Access Layer: The access layer includes client web browsers and applications. Browser-based users can access Developer Hub web front and its contents. Users can also view UEAPI events pages and perform searches on event data for a given time, location, or a topic. Application-based users can invoke the UEAPI web services using HTTP-JSON protocols.
2. Web Server Layer: This layer hosts Apache Web servers. To follow USPTO EA standards, Apache is configured as the web server in front of the JBoss EAP server. The Web Server layer serves two purposes—presenting Developer Hub’s static and dynamic content, and receiving and responding to UEAPI web services calls (HTTP Get/Post messages). The Apache web server routes the UEAPI web services calls to JBoss EAP, which hosts the UEAPI JAX Jersey RS RESTful Web services. JBoss EWS is the server and uses AWS Elastic Load Balancing (ELB) for load balancing applications.

3. Application Server Layer: This layer uses JBOSS EAP 6 JEE server to host UEAPI JAX-RS Jersey RESTful Web services and UEAPI backend services such as user authentication, email subscription / notification, ETL process and data synchronization. The JBOSS EAP servers are configured in cluster; if a server goes down, subsequent user requests can be forwarded to a different server.
4. Search Layer: To be defined by Release 3 of the project.
5. Data Layer: The Data Layer is responsible for providing access to the data from various sources, such as Drupal Relational Database (RDS), UEAPI Relational Database (RDS), Unstructured Events Data (AWS S3).
6. Infrastructure Layer: This layer provides user registration, authentication and authorization using MyUSPTO.

The Developer Hub Assignment Search (DH-AS) system indexes patent assignment records and allows them to be searchable by the public. To accomplish this, the system writes the internal records as files and transfers them to a receiving file system. A process monitors this file system and sends the records to the search system for indexing. Once complete with indexing, the whole file is transferred to another file system. If any errors occur, a third file system receives the file.

i) How information is transmitted to and from the system

BDR: Information is transmitted through batches, service calls, and user entry (BDR-TQR feature). All transmissions and retrieval of information are performed within the USPTO network and do not exceed the internal network boundary.

The BDR application employs a multilayered design approach. This approach gives modularity to the system. The following sections explain in high level, how each layer is comprised. The design principle of the BDR aims to have a tiered approach to the application. This way every component of the ecosystem is more easily understood and viewed independently. In this platform, there is ingestion, where the data is ingested from existing software resources. There is the compute tier, where the data is loaded, compared for public versus private status, and analyzed according to data science principles. There is the analysis tier, where data scientists combine the real world problem solving techniques from Patent Examiners with the formulae and hypothesis of the Data Science field. The Visualization tier that provides the users with a place to view the analysis and the underlying data that helps to create it. Finally, in the storage tier, the system retains raw, merged and transformed data, distinguishes between public and private Patent applications, and segregates them.

Developer Hub (DH):

The DH system provides USPTO public data (such as patents, trademarks, and events data) via a set of Web Services APIs for the consumption of the developer community. These APIs will be developed and maintained by various divisions within USPTO and will be accessible through a USPTO web UI named Developer Hub, or Davent Hub (DH) System Name.

The system provides access to USPTO public content through the use of APIs (application programming interface). It has been determined that DH does not process PII/BII information, and it is categorized as a low risk system. The DH web application is deployed on the Amazon Web Services (AWS) Cloud platform. Users include: General Public, System Development Staff, Tableau Public Users, EC2 Server Accounts, Drupal Admin User via RBAC, and System Administrators.

Developer Hub Assignment Search (DH-AS):

DH-AS is responsible for indexing patent and trademark assignment records, which allows them to be searched by the public. To accomplish this, the internal records are written as files and transferred from AHD to a receiving file system. DH-AS is hosted on an AWS Public Cloud using the IaaS Service Model. It has been determined that DH-AS does not process PII/BII information, and it is categorized as a low risk system. The DH web application is deployed on the Amazon Web Services (AWS) Cloud platform. Users include: PTONet internal users - Assignment Historical Database (AHD), Assignment Services Branch, USPTO personnel such as patent examiners and support staff, Public Search Facilities staff members, and SOLR administrators.

AS provides public access via Amazon’s Web Service Cloud the capability for external users of the USPTO as well as public users in the USPTO public search rooms (with access to the Internet) to query issued patent or published application patent assignment data and/or pending or registered trademark assignment data. The AS web application is deployed to the middleware environment running under Apache web servers and is available to external customers/users of the USPTO (outside of PTONet) via the Internet.

Questionnaire:

1. Status of the Information System

1a. What is the status of this information system?

- This is a new information system. *Continue to answer questions and complete certification.*
- This is an existing information system with changes that create new privacy risks. *Complete chart below, continue to answer questions, and complete certification.*

| Changes That Create New Privacy Risks (CTCNPR) | | | |
|---|--------------------------|-------------------------|--------------------------|
| a. Conversions | <input type="checkbox"/> | d. Significant Merging | <input type="checkbox"/> |
| | | g. New Interagency Uses | <input type="checkbox"/> |

| | | | | | |
|---|--------------------------|-----------------------|--------------------------|------------------------------------|-------------------------------------|
| b. Anonymous to Non-Anonymous | <input type="checkbox"/> | e. New Public Access | <input type="checkbox"/> | h. Internal Flow or Collection | <input checked="" type="checkbox"/> |
| c. Significant System Management Changes | <input type="checkbox"/> | f. Commercial Sources | <input type="checkbox"/> | i. Alteration in Character of Data | <input type="checkbox"/> |
| j. Other changes that create new privacy risks (specify): | | | | | |

- This is an existing information system in which changes do not create new privacy risks, and there is not a SAOP approved Privacy Impact Assessment. *Continue to answer questions and complete certification.*
- This is an existing information system in which changes do not create new privacy risks, and there is a SAOP approved Privacy Impact Assessment (version 01-2015 or 01-2017). *Continue to answer questions and complete certification.*
- This is an existing information system in which changes do not create new privacy risks, and there is a SAOP approved Privacy Impact Assessment (version 01-2019 or later). *Skip questions and complete certification.*

1b. Has an IT Compliance in Acquisitions Checklist been completed with the appropriate signatures?

- Yes. This is a new information system.
- Yes. This is an existing information system for which an amended contract is needed.
- No. The IT Compliance in Acquisitions Checklist is not required for the acquisition of equipment for specialized Research and Development or scientific purposes that are not a National Security System.
- No. This is not a new information system.

2. Is the IT system or its information used to support any activity which may raise privacy concerns?

NIST Special Publication 800-53 Revision 4, Appendix J, states “Organizations may also engage in activities that do not involve the collection and use of PII, but may nevertheless raise privacy concerns and associated risk. The privacy controls are equally applicable to those activities and can be used to analyze the privacy risk and mitigate such risk when necessary.” Examples include, but are not limited to, audio recordings, video surveillance, building entry readers, and electronic purchase transactions.

- Yes. *(Check all that apply.)*

| | | | |
|--------------------|--------------------------|----------------------------------|--------------------------|
| Activities | | | |
| Audio recordings | <input type="checkbox"/> | Building entry readers | <input type="checkbox"/> |
| Video surveillance | <input type="checkbox"/> | Electronic purchase transactions | <input type="checkbox"/> |
| Other (specify): | | | |

- No.

3. Does the IT system collect, maintain, or disseminate business identifiable information (BII)?

As per DOC Privacy Policy: "For the purpose of this policy, business identifiable information consists of (a) information that is defined in the Freedom of Information Act (FOIA) as "trade secrets and commercial or financial information obtained from a person [that is] privileged or confidential." (5 U.S.C.552(b)(4)). This information is exempt from automatic release under the (b)(4) FOIA exemption. "Commercial" is not confined to records that reveal basic commercial operations" but includes any records [or information] in which the submitter has a commercial interest" and can include information submitted by a nonprofit entity, or (b) commercial or other information that, although it may not be exempt from release under FOIA, is exempt from disclosure by law (e.g., 13 U.S.C.)."

Yes, the IT system collects, maintains, or disseminates BII.

No, this IT system does not collect any BII.

4. Personally Identifiable Information (PII)

4a. Does the IT system collect, maintain, or disseminate PII?

As per OMB 17-12: "The term PII refers to information that can be used to distinguish or trace an individual's identity either alone or when combined with other information that is linked or linkable to a specific individual."

Yes, the IT system collects, maintains, or disseminates PII about: *(Check all that apply.)*

- DOC employees
- Contractors working on behalf of DOC
- Other Federal Government personnel
- Members of the public

No, this IT system does not collect any PII.

If the answer is "yes" to question 4a, please respond to the following questions.

4b. Does the IT system collect, maintain, or disseminate Social Security numbers (SSNs), including truncated form?

Yes, the IT system collects, maintains, or disseminates SSNs, including truncated form.

Provide an explanation for the business need requiring the collection of SSNs, including truncated form.

Provide the legal authority which permits the collection of SSNs, including truncated form.

- No, the IT system does not collect, maintain, or disseminate SSNs, including truncated form.

4c. Does the IT system collect, maintain, or disseminate PII other than user ID?

- Yes, the IT system collects, maintains, or disseminates PII other than user ID.
- No, the user ID is the only PII collected, maintained, or disseminated by the IT system.

4d. Will the purpose for which the PII is collected, stored, used, processed, disclosed, or disseminated (context of use) cause the assignment of a higher PII confidentiality impact level?

Examples of context of use include, but are not limited to, law enforcement investigations, administration of benefits, contagious disease treatments, etc.


- Yes, the context of use will cause the assignment of a higher PII confidentiality impact level.
- No, the context of use will not cause the assignment of a higher PII confidentiality impact level.

If any of the answers to questions 2, 3, 4b, 4c, and/or 4d are “Yes,” a Privacy Impact Assessment (PIA) must be completed for the IT system. This PTA and the SAOP approved PIA must be a part of the IT system’s Assessment and Authorization Package.

CERTIFICATION

I certify the criteria implied by one or more of the questions above **apply** to the Open Data-Big Data Master System (OD-BD MS) and as a consequence of this applicability, I will perform and document a PIA for this IT system.

I certify the criteria implied by the questions above **do not apply** to the Open Data-Big Data Master System (OD-BD MS) and as a consequence of this non-applicability, a PIA for this IT system is not necessary.

| | |
|---|--|
| <p>System Owner Name: Scott Beliveau Office: OCTO-EAAB Phone: (571) 272-7343 Email: Scott.Beliveau@uspto.gov</p> <p>Signature:  Digitally signed by Users, Beliveau, Scott Date: 2021.05.05 08:06:47 -04'00'</p> <p>Date signed: _____</p> | <p>Chief Information Security Officer Name: Don Watson Office: Office of the Chief Information Officer (OCIO) Phone: (571) 272-8130 Email: Don.Watson@uspto.gov</p> <p>Signature: <u>DON R Watson</u> Digitally signed by DON R Watson Date: 2021.05.12 07:24:15 -04'00'</p> <p>Date signed: _____</p> |
| <p>Privacy Act Officer Name: John Heaton Office: Office of General Law (O/GL) Phone: (571) 270-7420 Email: Ricou.Heaton@upsto.gov</p> <p>Signature: <u>Users, Heaton, John (Ricou)</u> Digitally signed by Users, Heaton, John (Ricou) Date: 2021.05.04 18:37:17 -04'00'</p> <p>Date signed: _____</p> | <p>Bureau Chief Privacy Officer and Authorizing Official Name: Henry J. Holcombe Office: Office of the Chief Information Officer (OCIO) Phone: (571) 272-9400 Email: Jamie.Holcombe@uspto.gov</p> <p>Signature: <u>Users, Holcombe, Henry</u> Digitally signed by Users, Holcombe, Henry Date: 2021.05.12 13:48:07 -04'00'</p> <p>Date signed: _____</p> |
| <p>Co-Authorizing Official Name: N/A Office: N/A Phone: N/A Email: N/A</p> <p>Signature: _____</p> <p>Date signed: _____</p> | |