



# Building IP Analytics on Lighthouse IP's Structured Global Data Feeds

Lighthouse IP. The world's most complete IP data collection



# Contents

<u>Introduction</u>	2
<u>The commercial opportunity in IP analytics</u>	3
<u>Lighthouse IP's global data capabilities</u>	5
<u>Technical enablers for developers</u>	8
<u>APIs and advanced metadata</u>	10
<u>Use case examples</u>	12
<u>Architectural recommendations</u>	14
<u>Conclusion: partnering for innovation</u>	15
About Lighthouse IP	17

# Introduction

Lighthouse IP offers the only end-to-end global IP dataset of patents, trademarks, and designs delivered in clean, scalable formats. Its unique data collection and infrastructure enable analytics providers to focus on insights – not data wrangling – and thus launch differentiated products faster. The commercial opportunities are vast: from competitive intelligence dashboards and portfolio-scoring tools for law firms and in-house counsel, to prosecution analytics for patent attorneys, and IP-based investment signals for financial analysts and government agencies.

Lighthouse IP's structured feeds (via bulk XML on S3) are enriched with normalized legal-status and owner metadata, harmonized across jurisdictions, and linked at the entity level. This backgrounder explains how analytics companies can leverage these enablers (including patent-valuation APIs and USPTO prosecution metadata) to jump-start products and gain a market edge.

# The commercial opportunity in IP analytics

IP-rich organizations are increasingly demanding advanced analytics. Law firms and in-house counsel want tools to optimize patent prosecution and monitor competitors. Such tools require reliable examiner correspondence data, which Lighthouse IP's feeds can supply. Technology companies and startups need competitive intelligence solutions that reveal rivals' filing trends, citation influence, and patent ownership changes. Financial firms (hedge funds, asset managers, VCs) can use IP as alternative data: spotting innovation trends, litigation risk, or licensing events to drive investment strategies. Lighthouse IP's "IP for Finance" narrative shows how complete IP signals (from filings to transactions) become alpha-generating analytics. Government and economic agencies likewise benefit by tracking national patent portfolios and technology landscapes. Across all these sectors, one major cost is sourcing and cleaning IP data. Lighthouse IP's ready-made global feeds let analytics companies skip that step: plug in the data and start building insights immediately.

Key IP analytics applications include:

- Competitive intelligence & tech scouting: track where competitors and startups are innovating (e.g. identify new tech entrants via patent applications); monitor trademark filings for brand expansions.
- Portfolio scoring & valuation: use patent counts, citations, family size, and external valuation metrics to grade patent portfolios. IP-BI's Emposis platform, for example, provides per-patent qualitative profiles and monetary-value ranges that can be added via API. This supports merger due diligence and licensing negotiations.
- Prosecution strategy & patent quality: analyze prosecution histories (office actions, examiner rejections, amendments) to optimize drafting and examiner negotiation.

- IP-based investment signals: feed patent and trademark events into quant models (e.g. patent filings in a field, large citation networks, or sudden patent reassignments) as trading signals.
- Risk & M&A analysis: identify IP litigation and licensing events (legal-status changes), or patent value spikes, to inform litigation risk and M&A target selection.

By ingesting Lighthouse IP's complete, pre-normalized datasets, analytics providers can prototype and deploy such products much faster.

For instance, a patent CI platform might rapidly integrate a unified global patent XML feed (instead of stitching USPTO, EPO, JPO, etc. together) to immediately support cross-border patent-landscape queries. Similarly, a brand-monitoring service could launch global trademark alerts as soon as data arrives, without building out country-by-country pipelines. In each case, Lighthouse IP's feeds remove data overhead so teams can focus on algorithms and UI.

# Lighthouse IP's global data capabilities

Most complete worldwide coverage

Lighthouse IP sources from every major patent and trademark office – even those only publishing on paper. Its patent repository spans 170+ jurisdictions with 167 million bibliographic records. Its trademark collection covers 207+ authorities/registers, totaling 125 million records. It also maintains industrial design filings from 99 countries. In practice, this means no major country or emerging market is missing. (As one case study notes, customers gain rapid, scalable access to all offices they need.)

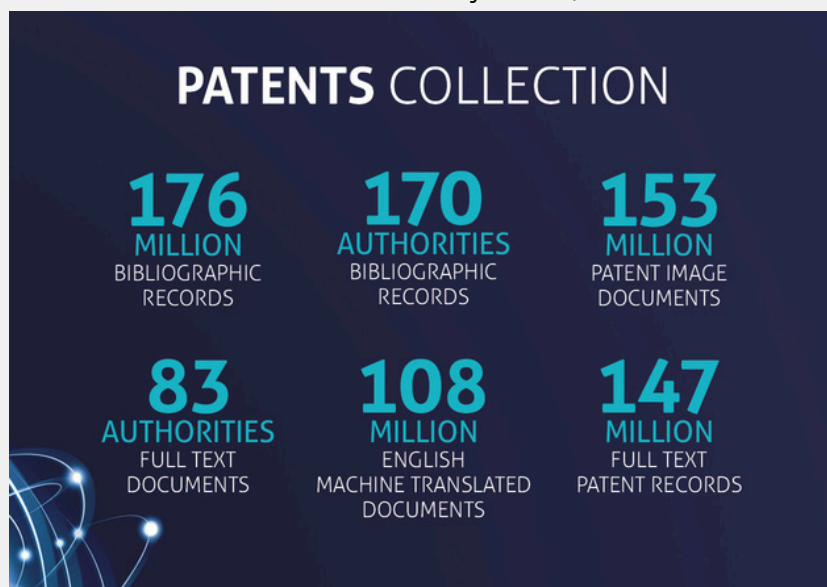


Figure 1: Lighthouse IP's patent collection: 176M+ records across 170 authorities

Behind this scale is the Diamond File concept: "insight into 170 authorities for bibliographic data and over 120 for up-to-date status information". Every published patent or trademark (with amendments, renewals, assignments, etc.) is captured and timestamped. For patents, Lighthouse provides full-text for approximately 81 countries and machine translations (English) from over 120 authorities, yielding 105 million translated documents. For trademarks, it delivers not only titles and classes but also images (112 million) and translations of goods/services lists where applicable. Crucially, Lighthouse IP's data is unified: it harmonizes naming conventions and status terms across all jurisdictions. This means analytics teams get a single-schema feed rather than a patchwork of heterogeneous records.

### Uniform structured feeds

Lighthouse IP delivers data in structured XML/JSON (based on WIPO standards) or CSV, directly to partners (typically via AWS S3). For example, all patent feeds use WIPO ST.36 XML (plus multi-page PDF images), while trademark data uses WIPO ST.66 XML.

Design data is also processed into a consistent proprietary format. Updates are pushed frequently (often weekly or daily) by authority. As one customer case notes, “the standardized provision of the data through Amazon S3 buckets improved timeliness, reliability and provided improved data security”.

Because everything is uniformly formatted, customers need no additional normalization. You can directly ingest these feeds into data lakes or databases without writing custom parsers for each country. (For instance, Lighthouse IP eliminated a customer’s need to maintain dozens of ETL routines by supplying a single standardized format.) The data pipeline is cloud-based and continuously maintained: any changes in source formats (say a PTO overhauling its XML) are handled by Lighthouse IP and delivered downstream without client effort.



Figure 2: Lighthouse IP’s trademark collection: 207+ authorities, 125M+ records

In practice, an analytics firm can spin up AWS infrastructure that simply pulls daily S3 dumps of Lighthouse IP's feeds. For instance, the firm might write an Airflow DAG or Kubernetes job that: (1) downloads patent XML from S3; (2) parses and stores bibliographic and claims data into a database; (3) similarly ingests trademark and design feeds; and (4) updates a unified "entity" table by matching owner names (Lighthouse IP provides translated, cleaned names). Such a pipeline might resemble an ELT for a data warehouse, with Lighthouse IP's pre-cleaned data skipping the usual "mapping and cleaning" stage.

#### Normalization and entity linkage

Lighthouse IP goes beyond raw documents by normalizing critical metadata. All legal status events (filings, grants, renewals, reassignments, oppositions, litigations, etc.) are explicitly captured and categorized. For example, Lighthouse IP's status feed is exhaustive: "if a patent was filed, pending, granted, litigated, or licensed, you'll know". Similarly, trademark updates include new registrations, amendments, renewals, and cancellations. These events are tagged as either new records or status changes, enabling analytics to track lifecycle trends.

On the owner side, Lighthouse IP uses proprietary translation algorithms to unify assignee/applicant names across languages and variants. For instance, a Japanese company name or Chinese applicant is machine-translated into a consistent romanization. Partners often use Lighthouse IP data to enrich their customer-company databases: Lighthouse IP itself notes that linking "company data [to] our extensive IP information" yields insights into each company's full IP strategy. In other words, if a corporation holds patents and trademarks across jurisdictions, Lighthouse IP's data connects them into the same entity tree. This entity-level linkage (essentially integrating all IP assets of a company group) accelerates analytics like portfolio aggregation, corporate benchmarking, and technology scouting.

# Technical enablers for developers

Lighthouse IP's platform is built on modern data infrastructure, enabling partners to focus on product logic instead of plumbing. Key technical features include:

- Bulk XML via S3: full historical (backfile) and daily feeds are delivered as compressed archives on Amazon S3 (optionally on HDD). Trademark backfiles and frontfiles, for example, "standard delivery for backfile and frontfile is per Amazon S3", with updates pushable via API or FTP. Patent data similarly can be bulk-downloaded. Teams can simply point AWS Glue, Spark, or EMR jobs at these S3 buckets to process the XML. This eliminates the need to individually crawl patent offices or perform FTP polling – Lighthouse IP handles all that.
- WIPO standard schemas: by conforming to WIPO XML (ST.36 for patents, ST.66 for trademarks), Lighthouse IP ensures that even though the sources are global, the data fields align uniformly. Your engineers won't have to reconcile "family" vs "publication" vs "kind codes" for each country. Instead, there is one schema per IP type. (For developers concerned about schema changes, Lighthouse IP promises to manage any source format shifts and supply updated loaders.)
- Normalized legal status fields: Lighthouse IP translates varied jurisdictional status codes into common categories (e.g. "patent pending", "patent granted", "renewal fee paid", "reassignment recorded"). This normalization is part of the feed, so analytics can directly filter or aggregate legal events. It also tracks deadlines: for example, it can flag when a patent is about to expire or a trademark renewal window closes. These normalized status flags support features like portfolio expiration calendars or maintenance cost projections.
- Entity resolution and owner metadata: all person and company names in assignee/applicant fields are disambiguated to the extent possible. For example, Lighthouse IP supplies an "owner ID" (when available) or consistently normalized name string for each record. (In one customer case study, Lighthouse IP highlights that translating assignees saves a Fortune 500 firm "tremendous effort".) This enables linking across datasets. E.g. you can join patent and trademark feeds on the normalized owner to see a company's complete IP footprint. Many analytics products need exactly this kind of entity-level mapping, and Lighthouse IP provides it out-of-the-box.

Integrated data types: unlike piecing together patent-only or trademark-only sources, Lighthouse IP offers full-scope coverage: patents and trademarks and designs. All three data types come in the same infrastructure. For instance, “One format vs. multiple silos” – with Lighthouse IP, “all patent data, trademark data, and design records merge seamlessly for analysis”. This means a brand-analytics tool could cross-reference trademark filings with patent filings of the same owner without wrestling with separate data vendors. In summary, Lighthouse IP’s platform is effectively a turnkey IP-data warehouse that feeds downstream analytics systems.

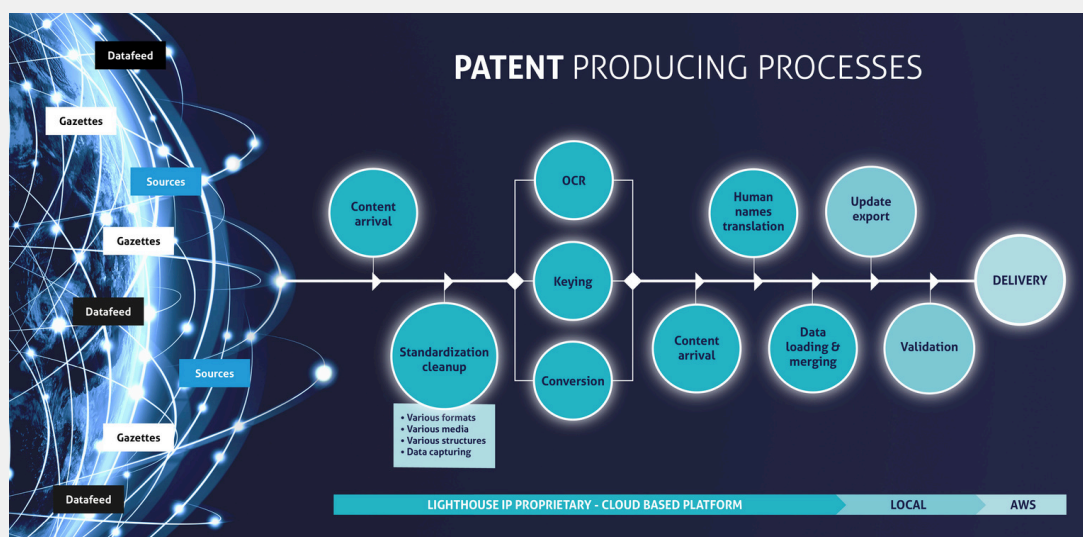


Figure 3: Global patent data ingestion pipeline. Raw patent publications (gazettes, feeds) are standardized via OCR, translation, and data cleaning into uniform XML for delivery.

In practice, a partner’s technical team might architect as follows: set up a cloud data pipeline (e.g. AWS EMR or Spark on Kubernetes) that ingests daily S3 shipments. The pipeline first validates that files match expected PTO formats, then parses the WIPO XML into database tables (or Spark DataFrames). Any necessary data-merging (e.g. linking family members into a single family ID) is done once, since Lighthouse IP assures up-to-date family structures. The pipeline would merge new data into an existing warehouse, while archiving raw feeds. Periodic checks (either Lighthouse IP-provided DDL or automated test scripts) ensure no fields change unexpectedly. Because Lighthouse IP “takes care of any format changes to connect as closely as we can into your data loading processes”, partners report minimal maintenance overhead.

# APIs and advanced metadata

Beyond bulk data, Lighthouse IP's ecosystem includes APIs and analytics services to enrich feeds:

- **IP-BI Patent Valuation API:**  
IP Business Information (IP-BI, a Lighthouse IP joint venture) offers patent valuation data. Their Emposis platform scores every patent family worldwide on multiple axes. For example, Emposis assigns each patent (or portfolio) a qualitative profile (metrics like market attractiveness, technical quality, legal strength) and a monetary value range. Analytics products can access this via data feeds or API to add an "estimated worth" to any patent. This is invaluable for portfolio scoring or M&A due-diligence features: an app could display a normalized patent-value score alongside the patent text. Lighthouse IP partners can license these valuations to overlay on their patent databases, giving end-users instant insight into which patents are likely most valuable.
- **USPTO File-Wrapper and Office Action APIs:**  
Lighthouse IP also recognizes the importance of prosecution data. Through industry-standard interfaces, developers can pull detailed USPTO information. For instance, the USPTO Open Data Portal provides an Office Action Text Retrieval API that returns full-text of examiner office actions. (Lighthouse IP data integration can consume this to append examiner communications to US applications.) Similarly, USPTO's Patent Application Oath/Signature (PA/AC) and Patent Assignment APIs can be leveraged for assignments/ownership history. In effect, an analytics team building prosecution tools can combine Lighthouse IP's global status feed (covering many jurisdictions) with USPTO's APIs for the deep US-specific history (file wrappers, petition decisions, appeal outcomes). This combined approach speeds development of modules like "patent prosecution dashboards" or "application risk flags".

- **Additional Linked Data:**  
Lighthouse IP's partnership announcements show they link to other databases. For example, a recent Lighthouse IP/IPDataLab collaboration enriches patents with FDA drug data, matching drug approvals to patent numbers. This highlights that Lighthouse IP's architecture can incorporate custom metadata via APIs. Analytics developers might likewise link third party business data or clinical-trial registries if needed, built atop the stable Lighthouse IP core.

In all, Lighthouse IP's APIs and feeds form a complete IP data backbone. With valuation scores and prosecution metadata readily available, an IP analytics product can deliver richer insights (e.g. "show patents with highest license revenue potential" or "prioritize patents with many office actions pending") without building these data pipelines from scratch.

# Use case examples

## Accelerated product development

Suppose an analytics startup wants to launch a patent portfolio valuation tool. Using Lighthouse IP, it could immediately access each company's global portfolio (bibliographic + status) and apply IP-BI valuation scores. Instead of months of sourcing data, the startup's engineers spend weeks building the scoring algorithms and UI. They might visualize, for example, how Company X's portfolio value changes year-over-year, all driven by Lighthouse's complete family and legal-event data.

## Competitive intelligence dashboard

A law-firm-focused CI product might ingest Lighthouse IP patents and trademarks to feed an interactive map of where each competitor is filing worldwide. With normalized assignee names, the tool easily aggregates filings by corporate group. Annual trends (e.g. "R&D ramp in Chinese patents for wireless") can be computed from the pre-cleaned data. For prosecution analytics, the product might call USPTO office-action APIs for US cases and merge that with Lighthouse IP's status feed for non-US cases. This could power a module, automatically highlighting which tech domains have unusually high rejection rates or long pendency.

## Brand expansion monitoring

A trademark analytics firm could utilize Lighthouse's TM feed (207+ authorities/registers, weekly updates) to offer next-day alerts when competitors register new classes or logos. Because Lighthouse IP includes machine-translated goods and services, the firm's AI can match brand terminology across languages. The product developer would spend negligible time on data integration, thanks to Lighthouse IP's unified trademark schema and S3 delivery.

## IP valuation and M&A advisory

Financial analysts often need the "most valuable patents" at a glance. By ingesting IP-BI valuations via Lighthouse IP, an analytics dashboard can instantly highlight a target's top-weighted patents or identify undervalued assets. Combined with Lighthouse IP's legal-status data (e.g. which patents have been heavily licensed), the result is a premium portfolio assessment tool. Such features can differentiate a product in crowded markets – and Lighthouse IP's data gives a head start over competitors who would otherwise struggle to gather this information.

### Government/research intelligence

A government R&D office could build a tech-forecasting tool on Lighthouse IP feeds, analyzing national patent output and key assignees. For example, it could cross-reference Lighthouse IP's patent and design filings to discover emerging industries. By leveraging the entity linkage, the tool could even correlate patents to parent companies or universities globally, supporting policy decisions on innovation hubs.

In each case, the first milestone of assembling and cleaning data is eliminated. Lighthouse IP's customers frequently report that having "a single data format for patents (and also for trademarks)" cuts onboarding costs dramatically. This lets analytics teams iterate on product features instead of debugging parsers.

# Architectural recommendations

Analytics products built on Lighthouse IP typically adopt a layered architecture. A common pattern:

- Data ingestion layer: automated jobs (on AWS EMR, Kubernetes, etc.) pull the daily S3 dumps. These jobs decompress XML files and stage them in a raw data store.
- ETL/processing layer: a Spark or Airflow pipeline validates schema, extracts relevant fields, and loads them into a database. Frequently, a star-schema is used (patents as fact table, with dimensions for dates, assignees, IPC classes, etc.). Because Lighthouse IP data is already prenormalized, this ETL step mainly focuses on performance tuning (e.g. partitioning by filing date).
- Analytics & service layer: this is where the innovation happens: machine learning models, search indices, and user interface backends operate on the curated data. For example, a patent clustering algorithm or a visualization engine would read from this layer.
- APIs & integration: above all, expose the intelligence to clients. Products often provide REST APIs or web dashboards. Many Lighthouse IP partners use BI tools like Tableau or PowerBI atop their warehouse (some may materialize aggregated views of the IP data for speed).

Throughout, Lighthouse IP's standardized XML/CSV feed acts as the single source of truth. Its consistency allows infrastructure-as-code (e.g. Terraform scripts) to be applied universally, rather than building country-specific ingestion code. The hosted S3 data can be versioned, so a team can backfill a year of patents, then switch to weekly mode seamlessly.

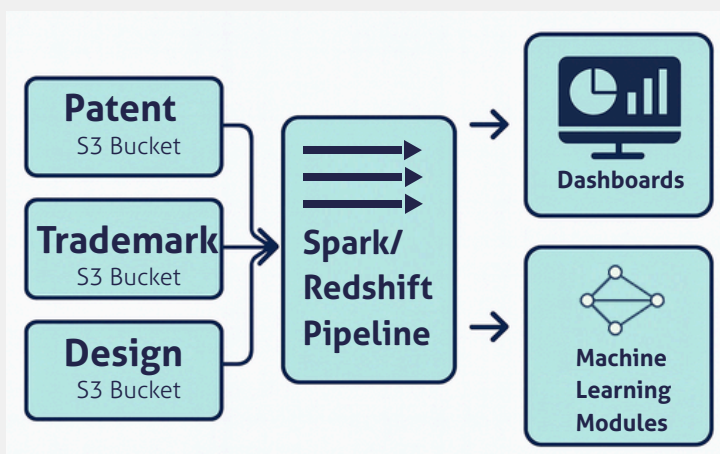


Figure 4: Lighthouse IP Data Architecture Flow

# Conclusion:

## partnering for innovation

Lighthouse IP's mission is to "unlock all published patent and trademark information". For analytics companies, this translates into a ready-made, enterprise-grade IP dataset that can be consumed via cloud services. By leveraging Lighthouse IP's structured global feeds, firms can accelerate time-to-market: skip the four months of data wrangling and dive right into building AI models, analytics dashboards, and decision-support tools.

Whether the goal is a new competitive intelligence platform, a sophisticated patent valuation service, or an AI-driven R&D scouting app, Lighthouse IP provides the reliable data foundation. As one Lighthouse IP client notes, "you no longer need to bother about getting data from many different sources in different formats". Instead, you can focus entirely on innovation on top of the data.

Next steps for interested partners

Contact Lighthouse IP to explore data licensing (whether patents, trademarks, designs, or combined feeds) and discover how IP-BI's valuation data and API services can plug into your architecture. With Lighthouse IP powering the data layer, analytics teams gain an unmatched competitive advantage – and the freedom to spark new product ideas without data constraints.

# About Lighthouse IP

Lighthouse IP is the world's leading provider of intellectual property content. We specialize in sourcing and creating unique data collections for patents, trademarks, and design information. Our processes cover every step: from acquiring original documents (in some countries still starting from paper) to delivering complete, uniformly formatted, and digitally accessible datasets.

With coverage of more than 170 countries and a global team of experts, we provide one of the most comprehensive and reliable bibliographic and legal IP data collections available.

Based in The Netherlands, Lighthouse IP was founded in 2006, with the same professionals who had previously successfully run Univentio. Key employees at Lighthouse IP have been in the industry for over 20 years. Our mission is to be the superior IP content provider in the world. Our data must be recognized as the most complete (largest backfile, number of content fields) and extensive (most countries) IP data collection, with the highest accuracy level. Our objective is to set the industry standard. As we source the data directly ourselves, we have several offices abroad, amongst others in Poland, China, Egypt, Indonesia, Thailand, the USA and Vietnam.