



CAS CUSTOM SERVICESSM
CASE STUDY

**UNRAVEL
THE POTENTIAL
OF INTERNAL
R&D DATA:
CURATE AND
CONNECT FOR
SEARCHABLE
INSIGHTS**

CAS



A division of the
American Chemical Society

Idea in brief

The knowledge you own has value, but what if you can't access it?

A large health-tech organization had decades of key research trapped in unstructured documents. As a result, relevant individuals struggled to locate, access, and apply this valuable information to their work. CAS helped this organization overcome data management challenges by leveraging the same expertise and technologies that make the world's scientific data available in CAS SciFinder[®]. As a result, internal experimental data is accessible, connected to the world's science, and enables better insights for data-driven decisions.



The challenge

When valuable data is inaccessible for research, discovery, and decision-making

Over decades, R&D teams amass large amounts of experimental results and findings in a variety of formats. This momentous amount of data that remains unknown, unused, and untapped across an organization can be known as 'dark data'.¹ It is estimated that 55% of an organization's data is 'dark'.² Generated by systems, devices, and interactions, the knowledge cannot be easily accessed due to its unstructured nature.

R&D data may sit in silos in several disconnected systems, making it difficult for individuals to access and utilize the information they need. This can lead to:

- Loss of legacy data gathered by individuals who have since left the organization.
- Limited ability to collaborate with individuals and teams across the organization.
- Inability to make connections between disparate datasets due to a lack of harmonization in scientific terminology.

90%

of industry leaders agree that every organization will need to extract value from dark data to be successful in the future.²

Splunk, The State of Data White Paper, 2019

80%

of the knowledge inside dark data is missed due to inefficient analysis.^{3,4}

IBM, 2018

When a knowledge management system is not fit for its purpose, it hides your own data from you and prevents you from making connections to external sources, leaving the potential for key insights to be missed.

While out-of-the-box solutions promise fast results, they are rarely able to accommodate the unique way organizations capture and store scientific data. The capabilities of any data platform are limited in their ability to deliver meaningful results without properly cleaning, normalizing, connecting, and enhancing the data first. Without specialized scientific and technical expertise, automated solutions will only go so far.

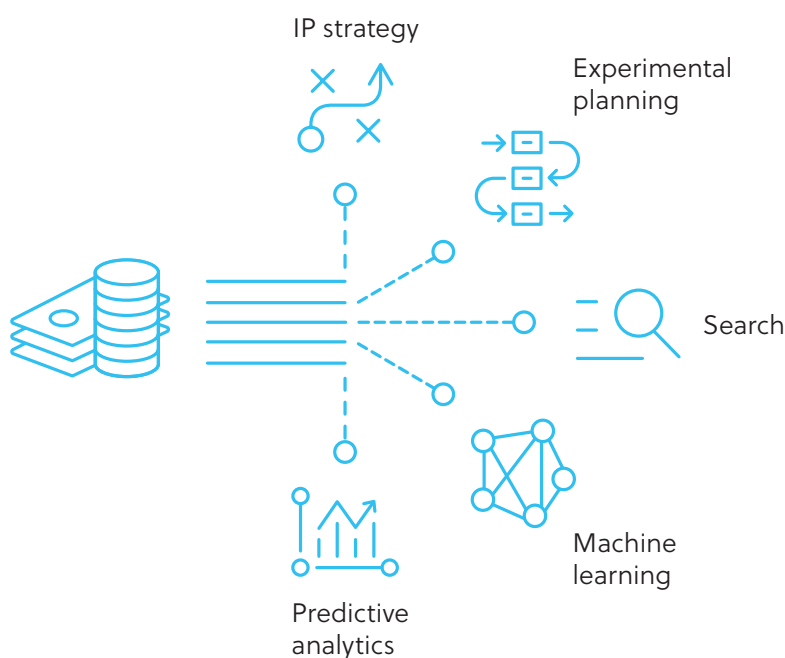


Figure 1. Unstructured, inaccessible data cannot be leveraged to drive future discoveries

Scientific complexity and connections: An ongoing challenge for in-house IT

The IT knowledge within R&D organizations is typically insufficient to deal with the complexities of scientific information; this is due to nuances with critical representations, nomenclature, and relationships. For instance, if the name and taxonomy of a drug has changed over time, a traditional database may not be able to make connections between varying information sources that have used inconsistent naming conventions. This may result in key information being omitted from searches. To access the information central to informing future discovery, technologies must be scientifically aware.

One organization's struggle

A large health-tech organization had decades of key research findings housed in several disconnected, legacy systems. There was no systematic way to search the documents; scientists had to know what they were looking for and manually sift through the files until they found it. The organization sought to better utilize decades of stored research through improved searching and retrieval of data, helping to drive more efficient innovation.

Stored proprietary knowledge was not being utilized often as veteran scientists retired and newly hired researchers were unaware of the available insights. Without access to these findings, many experiments were being unnecessarily repeated across R&D groups which cost considerable time and resources. Without the ability to access and utilize the knowledge gained by other researchers across the organization, innovation was hindered, costs increased, and schedules were pushed out.

The lack of organization, limited search capabilities, and an insufficient data model meant that valuable information remained hidden in their systems. This rendered it useless for future research, discovery, and decision-making.

"Having clear visibility to what data is already available is critically important. You need to be able to harvest the available information to reduce workload and speed up processes."

Director, R&D



The solution

Enabling access to your data through a custom solution

To enable the proprietary data to be readily leveraged, CAS built the client a custom knowledge management system to curate and connect their internal R&D data. This enabled the organization to search the full text of their documents, connect similar concepts and substances, and search by concepts tailored to their R&D efforts, such as indicating if a substance is a natural product for the health-tech organization.

This solution was made possible by utilizing the same process CAS uses to curate the world's published science that powers leading scientific information solution, CAS SciFinder[®].

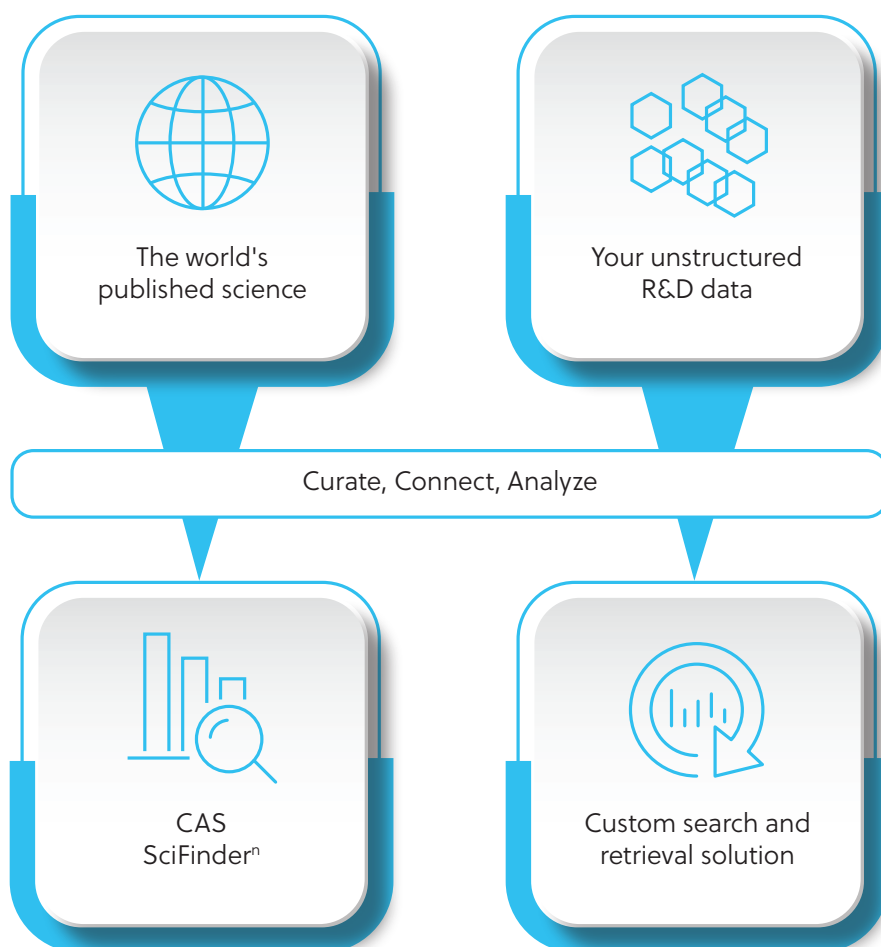


Figure 2. CAS reveals the full potential of your digital assets with similar processes that it uses to curate the world's published science

The CAS curation process utilizes the expertise of scientists alongside specialized technologies and approaches, such as:

- Tailored data governance and data management frameworks to ensure data is consistent and trustworthy
- Specialized lexicon, ontologies, and taxonomies to alleviate gaps from information nuances
- Machine learning algorithms, trained with CAS comprehensive data, that provide an accuracy edge to extract and connect data
- Manual curation for delivery of precise, high-quality data
- Optical character recognition (OCR) to digitize content
- Named-entity recognition (NER) & natural language processing (NLP) to identify substances and concepts

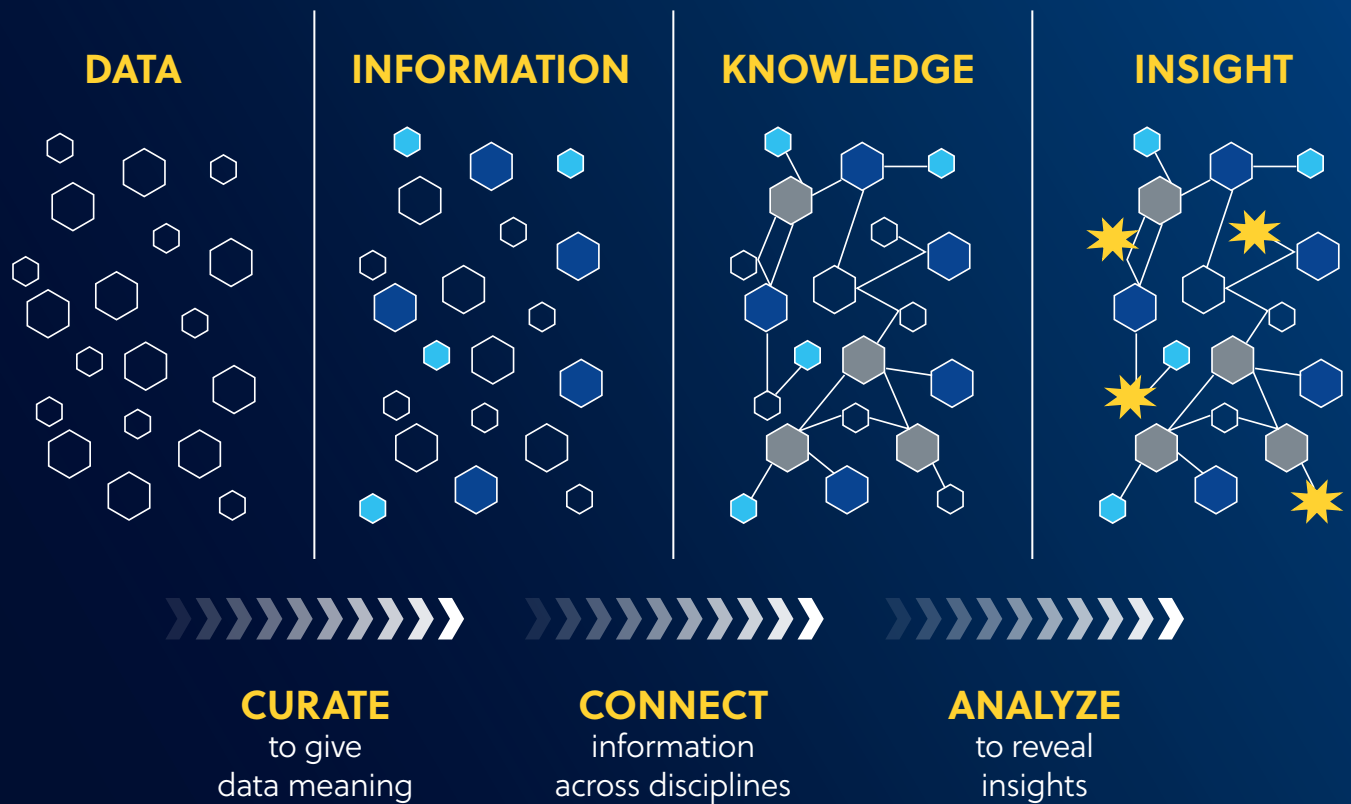


Figure 3. Data is valuable only when it is transformed into insight. Unlock the potential of your R&D data through capturing insights to help make data-driven decisions



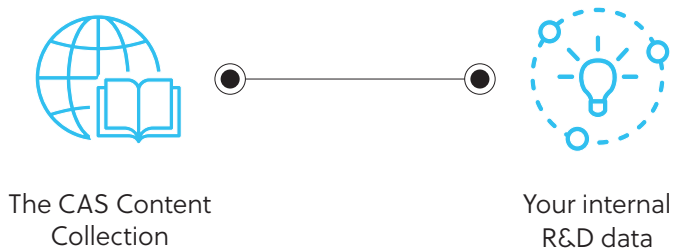
A scientifically aware search for meaningful results

To access the R&D information needed to inform decision-making, knowledge management solutions from CAS move beyond standard keyword search and enable scientific context. Information extracted from documents is associated with key concepts in chemistry, material sciences, life sciences, and custom criteria based on a client's innovation needs. A researcher can immediately find where a concept, such as bioactivity, polymers, or specific diseases, is referenced in the document.

To overcome the challenge of variations in synonyms amongst substances, scientifically aware searching enables users to take a trivial name and associate it with a structure not mentioned in the document. Refined lexicon, ontologies, and taxonomies allow for matching on synonyms (trade names, systematic names, internal names) or even matching on classes of substances. This sophisticated feature is far beyond the capabilities of any out-of-the-box solution.

To yield high-quality, impactful results, CAS utilizes sci-tech expertise paired with specialized tools and algorithms. Where human intervention is required, CAS can systematically route documents to trained domain-specific curation experts for conceptualization accuracy that only a scientist can provide.

Our proprietary substance comparison technologies use CAS Registry Numbers®, the worldwide standard for identification of substances, to disambiguate substance data. CAS Registry Numbers also connect internal data to the world's published science; this includes alerts of publications related to a user's internal research pushed to them in near real time.



Associated data

- CAS Registry Number
- Molecular properties
- Chemical structure
- Chemical names
- Regulatory information
- Reaction schemes
- Experimental procedures
- Additional data fields...

Figure 4. Scientifically aware search is enabled by connecting your organization's internal data with the world's published science.

Linking insights for internal collaboration and connections to the world's science

Knowledge graphs were constructed based on the available data enabling the organization to link information on substances, scientific information, and concepts, highlighting connections between the data across the organization. By utilizing knowledge graphs, researchers can find others in the organization working on similar areas of science, creating collaboration opportunities.

Data in internal documents is further enhanced when connected with the world's published science. Our unique data model securely links your propriety documents to the CAS Content Collection™ or custom-curated data relevant

to your industry, enabling you to identify and fill gaps in your data. Substance information becomes more robust when connected to data such as the CAS Registry Number, property data, and related research.

Referenced areas of science can be analyzed to determine trends in patents or publications, enabling you to understand who is publishing in your field and identify collaborators or key competitors.

Linking of related documents

Ion solvation scenario in an aqueous solution mixture of counteracting osmolytes: Urea and triethylamine-N-oxide (TMAO)

DOI: 10.1016/j.molliq.2019.111467
Pub Date: 2019-11-01

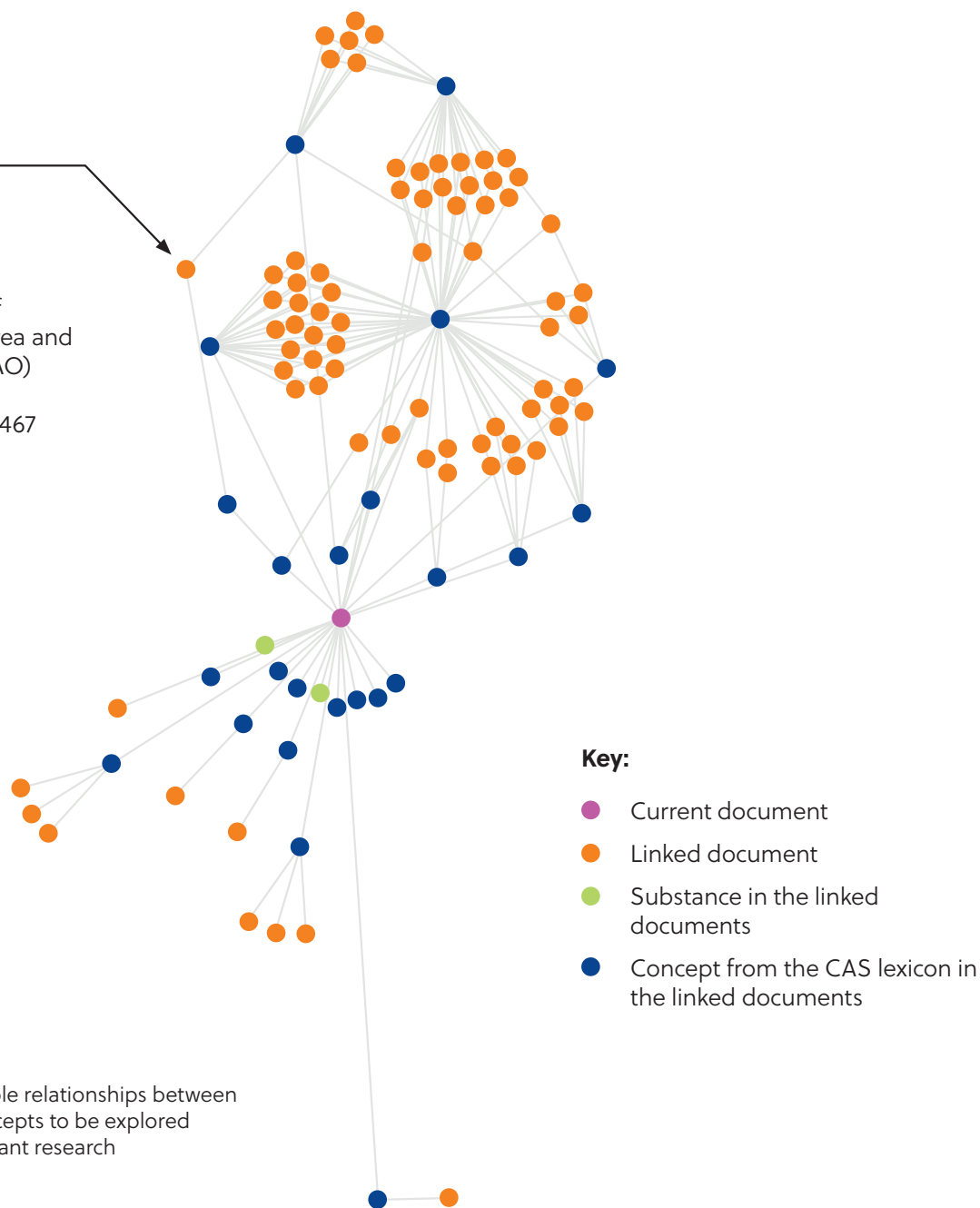


Figure 5. Knowledge graphs enable relationships between documents, substances, and concepts to be explored visually, connecting to other relevant research



Examples of a customized solution for searchable insights

isocyanate filtrate

Results for "isocyanate filtrate"

459 Results

Items per page: 10 1 - 10 of 459

1
Pauling lab notes - August 3rd, 1938

...with shaling at ou00b0C. solid dissolves slowly, deep yellow color of sol- (whereas isocyanate remains undissolved in n. NOOH alone). Bulky yellow pp. formed during reaction Faint...
...acidified - yellow pp. Filtered RP C. nearly white 10g 2 PM Aug. T Filtrate evaporated to 20cc. Nacl opted Discarded Properties Imol in SOH, isopropyl other, Ang...

Compare Original

2
Production of carbohydrate building blocks from red seaweed polysaccharides. Efficient conversion of galactans into C-glycosyl aldehydes

...trichloroacetyl carbamate derivative of 15 which was generated in the NMR tube by using trichloroacetyl isocyanate as a derivatizing agent (see Experimental Section). Moreover, the TLC profile of the... /-PrOH (90 mL), and then filtered through a glass-sintered filter. The filtrate was concentrated and coevaporated with toluene three times to give a brown-yellow solid (...
...suspended in MeOH (100 mL), and filtered through a glass-sintered filter. The filtrate was concentrated to afford a crude material (5.01 g) containing the dithioacetal derivative... (-PrOH (90 mL), and then filtered through a glass-sintered filter. The filtrate was concentrated and coevaporated with toluene three times to give a brown-yellow solid (... , 576-588 This Journal is u00a9
The Royal Society of Chemistry 2009 filter. The filtrate was concentrated to afford a crude material (5.00 g) containing the dithioacetal derivatives... 4-trichloroacetyl carbamate derivative of 15 was generated in the NMR tube by adding trichloroacetyl isocyanate (5 pL) to a CDCl3 solution of 15 (10 mg). 2.3... and acetone (100 mL), and then filtered through a glass-sintered filter. The filtrate was concentrated to afford a residue that was diluted with H2O (30 mL) and...

Compare Original

3
Pauling lab notes - January 13th, 1936

...yellow. PbloH) (ppted) added to remove sulfides. Turned brown. Filtered hot. yellow filtrate on cooling deposits pale yellow solid. Filtered Pp.light yellow, fietrate light yellow microscopic...
...replaced). Let cool over noon hour. Heated, filtered (greenish residue) fellow-orange filtrate. x lizing on cooling PbloHt added, heated filtered. Dork brown residue (PGS)!... yellow filtrate crystallizing atonce Jan 15, 1936 Certified, washed with water, centr., alcohol, centrifunged,...

Compare Original

4
Brønsted Acid-Catalyzed Highly Stereoselective Arene-Ynamide Cyclizations. A Novel Keteniminium Pictet-Spengler Cyclization in Total Syntheses of (±)-Desbromoarborescines A and C

...reaction mixture was diluted with CFCl3 and filtered through a small bed of Celite. The filtrate was concentrated in vacuo and the residue was purified using silica gel flash column chromatography. ... the reaction, the reaction mixture was filtered through a small bed of Celite and the filtrate was concentrated in vacuo. The residue was purified by silica gel flash

Diverse formats of reports

Handwriting

Figure 6. Search the full text of internal research documents that have been digitized and curated, including hand-written notes, tables, and a variety of data formats

CAS Indicative Examples isocyanate filtrate

← Return to Home

Found in This Document

Chemistry Materials Science Life Science Detection

diethyl 2-bromo-2-methylmalonate — (29263-94-3)

bromomalonate — (18917-18-5)

Knowledge Graph 2D 3D

JACS
JOURNAL OF THE AMERICAN CHEMICAL SOCIETY

Communication
pubs.acs.org/JACS

Enantioselective Organocatalytic Alkylation of Aldehydes and Enals Driven by the Direct Photoexcitation of Enamines

Mattia Silvi,[†] Elena Arceo,[†] Igor D. Jurberg,[†] Carlo Cassani,[†] and Paolo Melchiorre^{*,†,‡}

[†]ICIQ—Institute of Chemical Research of Catalonia, Av. Països Catalans 16, 43007 Tarragona, Spain
[‡]ICREA—Catalan Institution for Research and Advanced Studies, Pg. Lluís Companys 23, 08010 Barcelona, Spain

Supporting Information

ABSTRACT: Disclosed herein is a photo-organocatalytic enantioselective α - and γ -alkylation of aldehydes and enals, respectively, with bromomalonates. The chemistry uses a commercially available aminocatalyst and occurs under illumination by a fluorescent light bulb in the absence of any external photoredox catalyst. Mechanistic investigations reveal the previously hidden ability of transiently excited enamines to directly reach an electronically excited state upon light absorption while successively triggering the formation of reactive radical species from the organic halides. At the same time, the ground state chiral enamines provide effective stereochemical induction for the enantioselective alkylation process.

Figure 1. Mechanisms that enamines can use to drive the photochemical generation of radicals: (a) by inducing ground state EDA complex formation; (b) acting as a photosensitizer upon direct photoexcitation. The enamine radical cation, resulting from the SET, is not shown. The gray circle represents the chiral organocatalyst scaffold.

Herein, we demonstrate that the photochemical activity of chiral enamines and their potential for light-induced radical generation is not limited to the formation of ground state EDA complexes. This study unveils the to date hidden ability of enamines, generated by condensation of aldehydes with a commercially available chiral secondary amine catalyst, to reach an electronically excited state upon simple light absorption and then to act as effective photosensitizers (Figure 1b). At the same time, such photochemical behavior conjugates with the enamine ability of stereoselectively intercepting the reactive radicals generated from the sensitization of organic halides. This novel photo-organocatalytic mechanism overcomes the need for an external photoredox catalyst. We used it to develop the enantioselective alkylation of aldehydes with bromomalonates **1**, a transformation that served as a benchmark for developing enantioselective dual photoredox–organocatalytic systems.^{24,47}

Identification of chemical substances

Figure 7. Substances, terms, and concepts are algorithmically extracted and embellished, linking key scientific data

References

1. *What is Dark Data?* https://www.splunk.com/en_us/data-insider/what-is-dark-data.html (accessed 2022-09-30).
2. *Splunk>. The State of Dark Data.* <https://www.splunk.com/pdfs/dark-data/the-state-of-dark-data-report.pdf> (accessed 2022-09-30).
3. *Dark Data Analytics.* <https://articles.xebia.com/dark-data-analytics> (accessed 2022-10-05).
4. *Illuminate Dark Data – Part 1.* <https://medium.com/ibm-data-ai/illuminate-dark-data-part-1-c3c2028f91af> (accessed 2022-10-05).



The impact

Time and cost efficiencies

With proper strategy and structure, research findings previously hidden in legacy technologies are now readily accessible to the R&D team. Before implementing the customized search and retrieval tool, one search for a specific data point within the unstructured files took a senior scientist approximately eight hours; now it takes minutes. This not only frees up resources for more specialized tasks, but it also creates cost savings by preventing unnecessary experimentation and repetition of work.

A competitive advantage

Data drives decision-making. With access to the full breadth of the organization's knowledge, teams can now leverage their valuable R&D data to accelerate the innovation of new health-tech solutions. Our customized solutions enable organizations to quickly respond to challenges that arise, from supply chain problems to the duration of reformulation. As a result, time to market can be reduced, ensuring your organization remains at the forefront of cutting-edge R&D.

Looking to unlock the potential of your data?

CAS Custom Services designs solutions so you can take advantage of the information you own to drive discovery. Reach out today at cas.org for an initial consultation.

One search took a senior scientist ~8 hours; now it takes minutes with scientifically aware search from CAS.

CAS is a leader in scientific information solutions, partnering with innovators around the world to accelerate scientific breakthroughs. CAS employs over 1,400 experts who curate, connect, and analyze scientific knowledge to reveal unseen connections. For over 100 years, scientists, patent professionals, and business leaders have relied on CAS solutions and expertise to provide the hindsight, insight, and foresight they need so they can build upon the learnings of the past to discover a better future. CAS is a division of the American Chemical Society.

Connect with us at cas.org

CAS



A division of the
American Chemical Society