



DATA RISK REPORT

Q1 2023

AUTONOMOUS DATA SECURITY

TABLE OF CONTENTS

EXECUTIVE SUMMARY
DATA RISK FINDINGS & TRENDS
EVALUATING RISK
DANGEROUS PATTERNS
OVERSHARING AND OTHER RISKY MIS-STEPS
REAL WORLD DATA RISK EXAMPLES
CONCLUSION

“

Concentric autonomously assigns data to one of over 250 categories. Over 85 of those categories are business-critical.

”

EXECUTIVE SUMMARY

Thus is the 2H 2022 edition of the semi annual Data Risk Report published by Concentric AI

Over 80% of an organization’s data is unstructured, meaning it’s embedded in the millions of financial reports, corporate strategies documents, source code files, and contracts created by CFOs, general managers, engineers, and lawyers every year. But to an IT security professional, unstructured data is still a shapeless lump of clay – unformed, unseen, and insecure. Most enterprises lack visibility into where their sensitive data is, much less where the risk is to the information from entitlements, sharing, permissions, activity etc.

Using advanced AI capabilities, Concentric processed 500 million unstructured data records files from companies in the technology, financial, energy and healthcare sectors. This report gives shape to the state of risk to unstructured data in the real-world by categorizing the data, evaluating business criticality, and accurately assessing risk. Accuracy is critical – it’s the difference between effective protection and alert fatigue caused by thousands of false positives.

All results in this report were based on live data analyzed by the Concentric Semantic Intelligence™ solution autonomously derived to reach our conclusions. Here are a few things we learned and how the data compares to 1H 2022

Concentric identified 250 biz critical categories (flat from last cycle)

- Nearly 32% of an organization’s unstructured data is business critical (meaning its distribution should be controlled)
- 90% of business-critical documents are shared outside the C-suite
- Over 15% of all business-critical files are at risk from oversharing, erroneous access permissions and inappropriate classification and so can be seen by internal or external users who should not have access

15TB	Unstructured data per avg enterprise
5TB	Biz critical data per avg enterprise
0.75TB or 402 data files per employee	At risk biz critical data per employee ((up from 310 in 1H 2022)

OVERALL

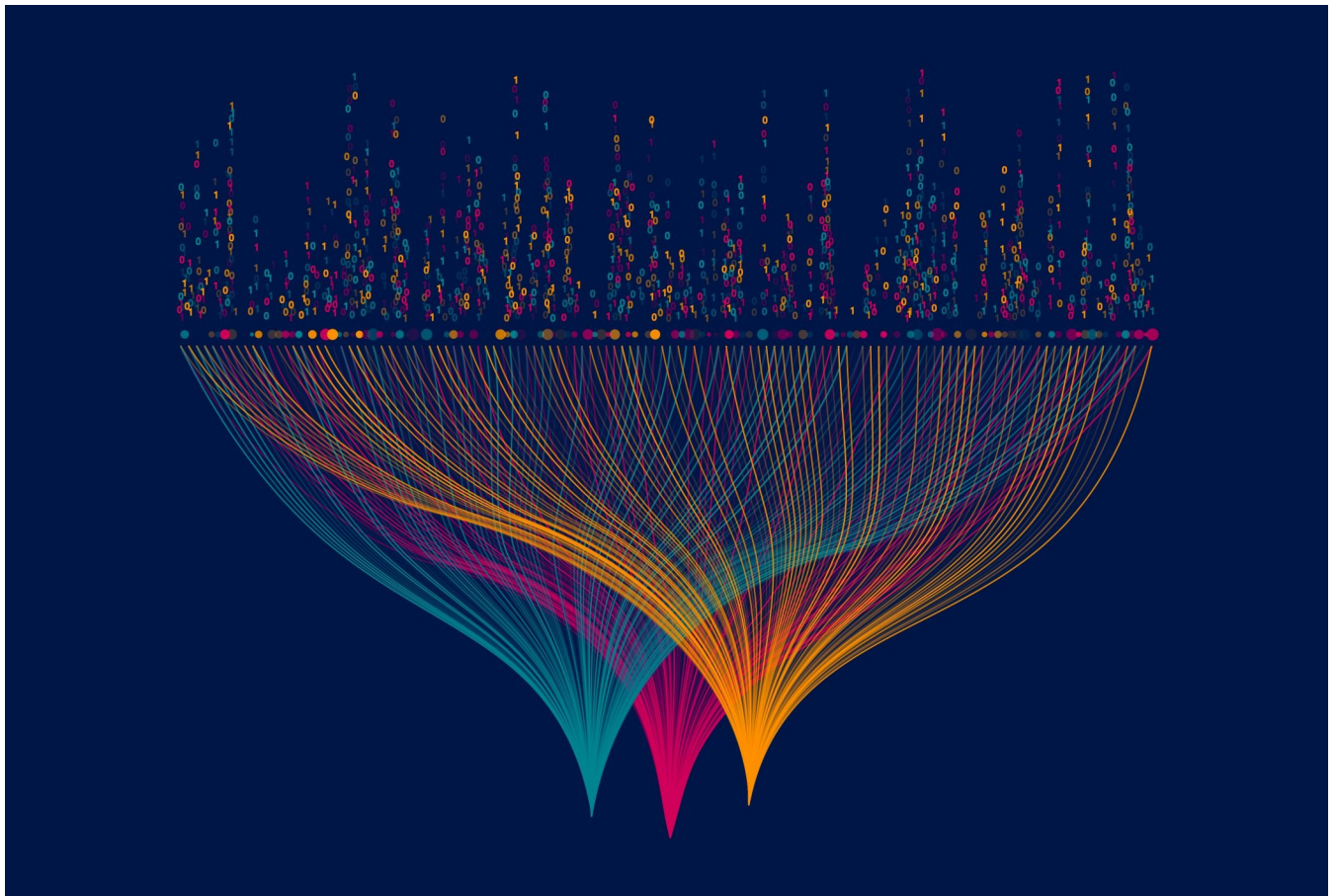
- Risk due to oversharing continues to trend up – a 12% HOH increase. Some of it is attributable to the increase in sample size but adjusting for that, there has been an increase in risk due to sharing. The data proves that in-spite of all the cyber security investments, data remains a vulnerable threat surface
- On average, each organization had 802 K data files at-risk due to oversharing (402 files per employee) up from 598K in 1H 2022 (310 data files per employee)
- Link based risky sharing is up to 100K documents per enterprise (from 81K documents per enterprise in 1H 2022)

	1H 2022	2H 2022	Delta
Total Unstructured data analyzed for this report	150TB	500TB	+200%
Biz Critical data files per enterprise	4.26M	5M	+17%
Data at risk from Oversharing	598K	802K	+34%
Data at risk per employee	310	402	+30%
Avg Employees per customer	1932	1995	

The Great Security Gap: Unstructured Data

Today, perimeter control and database protection products get the lion's share of security spend. Firewalls, access control frameworks and cloud access security brokers (along with many other security solutions) are large, established product categories. Enterprises have options. But the options to protect unstructured data aren't nearly as focused or effective. It's

easy to see why scanning PBs of unstructured data to accurately identify those that are both business critical AND inappropriately shared is no small feat. This, then, is the crux of the unstructured data security problem. Out of the 15.2 million files an average enterprise has, how can we know which files are overshared without overwhelming IT teams with false positives?



Taking Shape

Unstructured data is diverse, both in form and content. Many files are mundane and represent no real threat if overshared or stolen. Others contain information critical to the business. So, the first – and perhaps most difficult – task is to determine which documents we should worry about.

Concentric AI used sophisticated deep learning techniques to categorize over 500 million files from companies in the technology, finance, energy, chemicals, university and healthcare sectors. We discovered that the average organization has over 251 different types of biz critical categories hidden in its unstructured data (grouped here for clarity)



Product

(representative categories include bills of materials, source code, design documents, and test plans) - overshared product files can result in a loss of intellectual property, increased product liability, customer anxiety, and strategic disclosure.



Human Resources

(representative file categories include offer letters, stock agreements, and consulting contracts) - oversharing HR files can harm employee satisfaction, reveal private information, and driver higher costs



Financial

(representative file categories include bookings, income, revenue forecasts, pricing documents, invoices, trading, and tax filings) - oversharing these files can result in insider trading violations, compliance liabilities, and loss of competitive advantage.



Sales

(representative file categories include requests for proposals, quotes, and customer strategies) - exposing sales files can result in lost business, strategic disclosure, and sales team dissatisfaction.



Legal

(representative file categories include non-disclosure agreements, contracts, and purchase agreements) - exposure of a sensitive legal file can expose the company to civil lawsuits, loss of favorable supplier terms, and other legal liabilities.



Partner

(representative file categories include mergers and acquisition documents and partner agreements) - losing partner files can damage partner relationships, sink acquisition initiatives, or encourage insider trading.

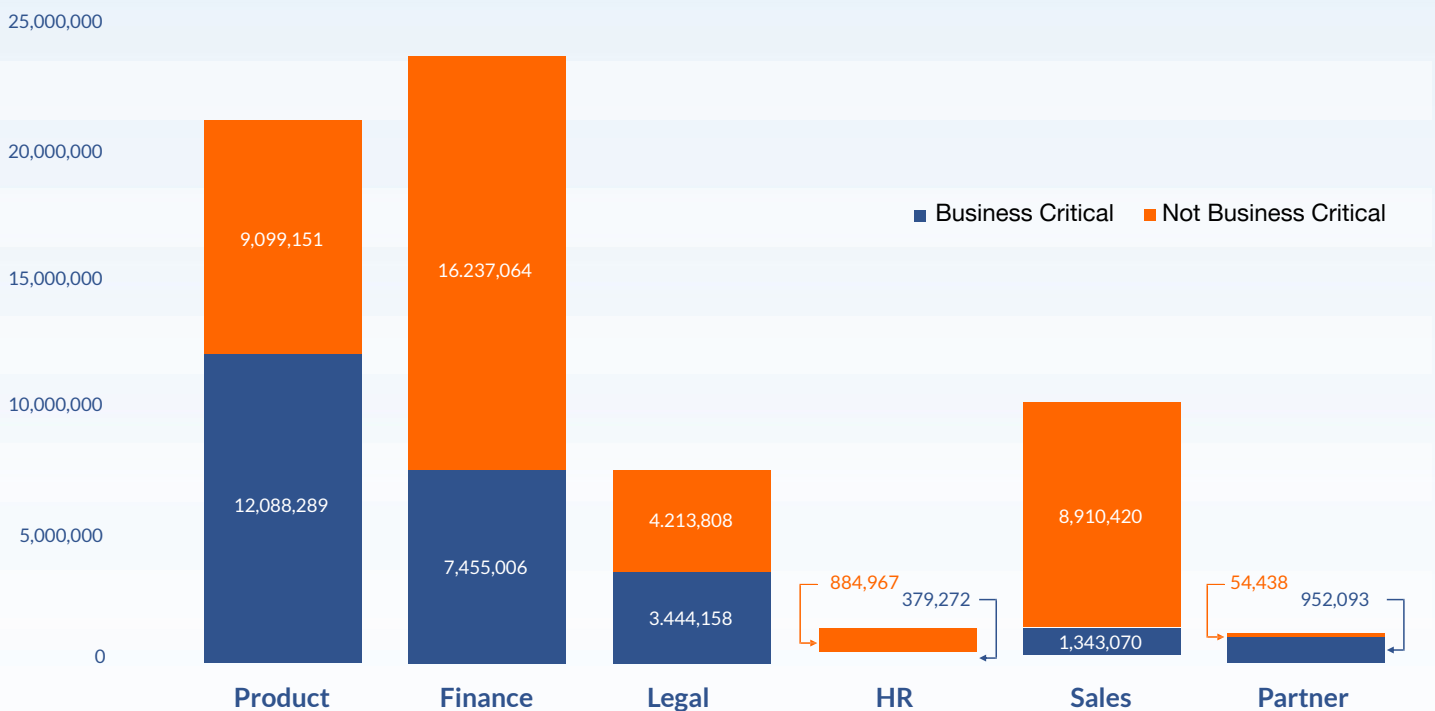
Seeking Meaning

Once categorized, Concentric AI evaluated each file for business criticality based on a variety of factors including , contextualized content, file ownership, document metadata, presence of personally identifiable information, and peer file comparisons. Business criticality is, of course, a vital piece of the puzzle. These are the files that must not be overshared.

HERE'S WHAT WE LEARNED:

- Nearly 32% of an organization's unstructured data is business critical (5M million files on average per organization)
- On average, each employee is responsible for 2506 business critical documents
- Financial data accounted for the leading share of business-critical documents (24%), followed by product files (22%) and sales files (10%) and legal documents (8%)

Unstructured Data by Category



Finding Answers

Assessing risk – even with a fully categorized set of files - is a deceptively complex task. Appropriate sharing depends on the meaning and function of the data record itself.

Sharing a contract with the legal team may be appropriate. Sharing it with the engineering team might not. It depends. Rigid rules- based risk assessments can be wildly inaccurate, especially when applied to policies governing intra-company data sharing.

To gain an accurate picture of risk, our analysis starts with the categories developed in the previous steps. We compare each document's security parameters to those of its peers. Using peer data configurations as a benchmark we can reliably identify oversharing – especially between employees at the same company or with external 3rd parties.

We check for:

- **Sharing with external users** – are similar documents shared with external users? If so, are they the same external users?
- **Sharing with groups** – do peer files allow similar group access?
- **Sharing with internal users** – is internal user sharing consistent? This is tough to spot without peer file analysis – and it's critical for security.
- **Misclassified confidential files** – has this document been properly classified? Document metadata, such as a “confidential” tag, is routinely used by other security solutions to enforce policy (e.g. a DLP solution uses a tag's setting to block a document from inappropriate access)
- **Misclassified files containing PII** – is this document marked to indicate it contains PII? Classifications for PII can also help a DLP solution fence in PII to maintain privacy and compliance
- **Wrong location**
- **Anonymous link sharing**
- **Sharing with personal email accounts**

COMMON SCENARIOS THAT INCREASE RISK

External user sharing mismatch	Biz critical documents shared inappropriately with external users
Group sharing mismatch	Sensitive data shared erroneously with groups
Internal user sharing mismatch	Sensitive data shared erroneously with internal users
Misclassified documents	Confidential documents that have been misclassified and can be accessed by the wrong personnel
Unclassified and PII	PII in docs that have not been classified
Wrong location	High value data in the wrong location
Personal email sharing	Data shared with personal emails
Risky Link sharing	Data shared via anonymous link sharing

We discovered some surprising results:

- 16% of an organization's business critical data is overshared
- On average, each organization had 802k files at-risk due to oversharing (402 files per employee) up from 598K in prior qtr (310 files per employee)
- Documents in the product and finance categories accounted for 41% of the total number of overshared documents
- 83% of the at-risk files were overshared with users or groups within the company (flat from prior quarter)
- 17% were overshared with external 3rd parties
- 87K business-critical files were erroneously classified and accessible by employees who should not have access to it

Oversharing per enterprise	1H 2022	2H 2022
Oversharing with external users	81,610	17,7945
Oversharing with internal users	209,539	278,195
Oversharing with internal groups	156,603	187,970
Oversharing due to misclassification	59,553	87,719
Oversharing due to wrong location	143,369	192,982
Oversharing with personal email accounts	11,028	52,632
Anonymous link sharing	19,851	30,075

DANGEROUS PATTERNS

After reviewing the data, we noted some patterns that seemed to reoccur across companies, regardless of sector or company size.

Near Duplicate files.

1 in 3 files we processed were identical or nearly identical. Near Duplicate files create multiple variant copies of sensitive information, often with different (and incorrect) file permissions, prohibited locations, or improper file classifications.

Shared with everyone

A shocking number of business-critical files were shared with everyone in the company. We found 160,078 such files

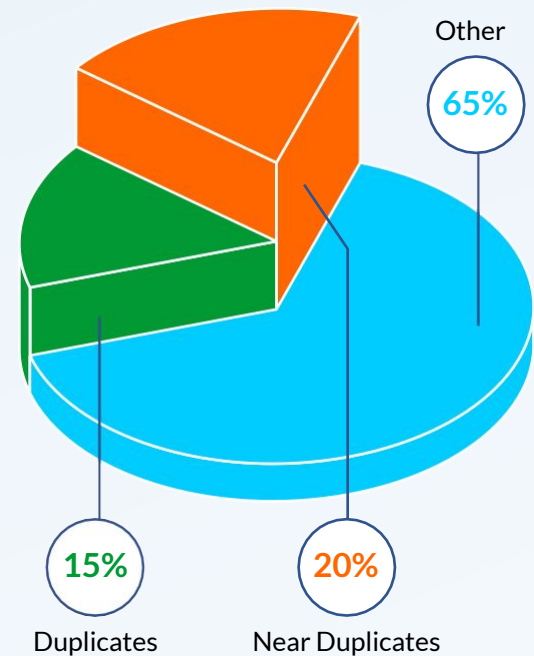
Internal or external oversharing

Of the 802K at-risk files, 83% were overshared with users or groups outside the proper team or department. These issues are nearly impossible to identify without peer data comparisons.

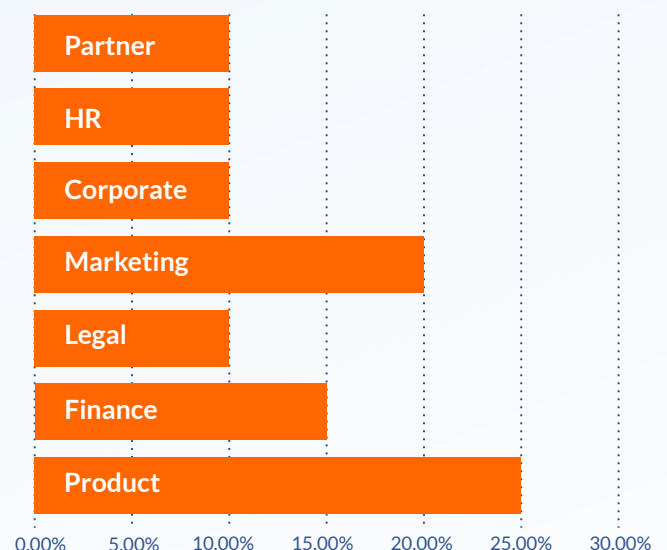
Personally identifiable information (PII)

PII is of interest due to privacy concerns and regulatory requirements. Increasingly, security teams use document metadata to flag PII and control document sharing and transfer. Nearly 25% of all documents containing unstructured data contained PII and were not marked appropriately.

Duplicates and Near Duplicates



PII as percentage across categories



Oversharing is a Modern Enterprise Reality

We've provided no shortage of statistics in this report. Statistics, sometimes, don't convey what's really happening on the ground. To help keep it real we offer a few specific incidents that show just how easily oversharing happens.

Risky sharing outside the company

Bob in Finance at an energy firm shared sensitive ITAR protected data with a friend

Risky sharing to personal email

The CFO at a high tech company needed access at home and she sent the confidential 2023 budget to her Gmail account

Not enough access controls on IP or PII data, ie. too permissive

Judy in HR at a financial services company ended up having access to highly-sensitive financial intellectual property

Sensitive data in the wrong location

At a healthcare firm, we found healthcare documents with PHI stored in Office365, Google Docs, and Dropbox when they wanted it securely stored only on AWS S3

Inappropriate classification

At a financial services firm, we found mortgage documents that were not classified correctly and consequently open to almost everyone in the IT department

Conclusion

Data sharing's transformational impact on businesses is indisputable. The productivity of pre-network paper-based communications pales compared to today's instantaneous electronic world. Hot new technology trends—like cloud computing and corporate digital transformations – all advance one fundamental goal: more, and more effective, data sharing.

But sharing has its dark side. The documents that used to fill physically secure filing cabinets can now be shared instantly, and with anyone. Businesses risk oversharing confidential sales strategies, need-to-know M&A plans, and sensitive personnel information with people who shouldn't have access. Autonomous and intelligent technologies are now able to find, categorize, and assess large bodies of unstructured data for better data security posture management and security.