



Migrate to the Cloud

THE HOW AND WHY OF MODERNIZING YOUR DATA WAREHOUSE



CHAMPION
GUIDES

What's inside:

- 3 What the market says
- 4 Part 1—Migration framework—approaches, strategies and requirements
- 6 Strategy requirements
- 8 Three approaches for implementing your migration
- 10 Part 2—Planning your migration
- 14 Part 3—Taking advantage of your new, cloud data warehouse
- 16 Conclusion
- 16 Learn more
- 17 About Snowflake



What the market says

According to Gartner, 92 percent of organizations are aware they have unmet data management demands in support of their analytics.¹ Seventy-five percent of data executives say they can't deliver useful data analytics effectively to the enterprise due to inflexible computing solutions². As a result, more than one-third of data professionals say their organizations are already using a cloud data warehouse³. But this doesn't mean they want a "cloud-washed" version of a legacy and inflexible on-premises data warehouse. Nearly all data professionals, 93 percent, see the unique benefits of a data warehouse built from the ground up for the cloud⁴.

Enterprises realize that legacy data warehouses can no longer deliver on their true purpose to organize data, enable rapid analysis and make insights available to all business users who need them. That's why they're moving away from traditional data warehouse solutions toward cloud solutions.

With some upfront planning and consideration, migrating your data analytics to the cloud is a process that can lead to big payoffs for your business and technology demands. In this eBook, we'll address your organization's data analytics needs with a roadmap for migrating your data warehouse to the cloud.

¹ Survey Analysis: New Data and New Analytics Are All Mythology Unless You Add Skills, Gartner.com, 9/18/17

²⁻⁴ Data Analytics: Beyond the Hype. A Survey of Data Professionals and Executives, Dimensional Research, 9/16



Migration framework— approaches, strategies and requirements

There are many reasons organizations choose to embrace cloud computing. But most organizations need a plan, something to grab onto and see what the future looks like. Not everyone is in the same place regarding their analytic capability or cloud maturity. Therefore, give careful consideration as to how fast and how much legacy code you want to move from your on-premise environment to a public cloud infrastructure.

THE FOUR MOST COMMON MIGRATION SCENARIOS

The type of migration you embark on will significantly influence your migration strategy. Here are four potential paths many organizations take to migrate their data analytics and data warehouse to the cloud:

1. OLTP for operational reporting and analytics

This is extremely common. Many organizations use OLTP (online transaction processing) systems, such as SQL Server, Oracle, or MySQL for basic reporting and analytics. While this might work as a short-term solution, the reporting needs of the business compete with the operational needs, overtaxing a fixed resource and slowing performance for both. A truly elastic cloud data warehouse eliminates this problem. It's pretty easy. As discussed later, take your existing transactional schema, which is usually in 3rd normal form, and move it, as is, to the cloud. This removes the reporting workload from the existing system and houses the data in a platform built for analytics. This eliminates the performance bottlenecks, and in some cases, gives your operational data store new life.

2. Appliance-based data warehouse

Many of the on-premise appliance vendors are sunsetting these legacy systems. More importantly, their customers want to escape the performance, cost and other limitations of these systems that can't address the ever-changing data analytics needs of the modern enterprise. Appliance-based data warehouses also require huge upfront costs typically in the form of capital expenditures (CAPEX).

On the flipside, a data warehouse built for the cloud is designed to grow with you from zero upfront costs thanks to a pay-as-you-go model, representing an operational expenditure (OPEX). This removes the guesswork of planning for your biggest day of consumption and then overpaying for an underutilized system for the other 364 days of the year. Similarly, if you need to expand your analytics unexpectedly during the year, you are hamstrung by a system that can't dynamically adjust to meet your needs. Another benefit is near, real-time access to data. Since on-premise appliances are a fixed resource, data warehousing teams create overnight load windows to make data available for the next morning. Today's cloud-built technologies allow you to segment workloads and load data 24/7 without impacting query processing, speeding the time-to-value of your data.

3. Data marts but no data warehouse

Most organizations suffer because a single source of truth is always out of reach. They have data siloed in many repositories. They may have tried to federate access across these repositories but quickly realized the cost to create and maintain that access wasn't feasible. They need a centralized repository to eliminate these barriers to getting all the insight from all their data. A cloud-built data warehouse becomes an obvious choice for data consolidation because it's ACID-compliant (transactionally consistent), can be partitioned/segmented logically without replicating, and can scale computing resources up and down, and on-demand. Whether you're a Kimball or Inmon fan, you need a platform that separates compute from storage and allows end users the greatest flexibility to access data sets using enterprise tools that leverage ANSI SQL.

4. Data lakes—data in, no insight out

Data lake initiatives have become a proverbial black hole: easy to get data in, complex but impossible to get data out. Many enterprises have realized that on-premise Hadoop infrastructures are costly and complex, and don't meet their analytics and concurrency requirements. Fortunately, there is a path forward. Leveraging a "layered" or "zones" approach is a great way for an enterprise to identify data sets they can comfortably move to the cloud. This method makes it easy to show the movement of data from on-premise to the cloud in a controlled and secure manner into cloud storage infrastructures such as Azure Blob Storage and AWS S3.



Strategy requirements

Migrations aren't much different than most IT projects, which means they usually begin with requirements. Defining requirements often cross multiple boundaries since a cloud migration strategy can be an executive-level decision. Without executive buy-in, your project will be limited in scope and be labeled as "shadow IT." This might work for some lines of business to get started. Eventually, everyone needs to be on the same page, from the architecture team to the security team to the chief data officer and even the CFO.

ORGANIZATIONAL STRATEGY

For many organizations, cloud is in their DNA and even written into their mission statements. For others that have been around awhile, they know they need to modernize but not at the expense of changing too quickly. The business benefits of the cloud are hard to ignore – more agility, lower costs, deeper analytics. But your cloud migration project should move at a pace consistent with your corporate

objectives. Determine if the project is the "tip of the spear", a way to get the company moving towards the cloud in the right direction. Or, part of a larger cloud initiative which would allow sharing best practices and technical resources. One of the critical success factors will be determining which data sets are ready to leave your data center first and which data sets will follow later to the cloud.

TECHNICAL STRATEGY

Every strategy starts with the same question: "What are we trying to do?" For technology people, this can be defined in the requirements. There is an old adage that says, "You can have it fast, good or cheap. Pick two." Using a combination of agile development, and taking advantage of what the cloud offers, that adage is more antiquated than accurate.

BUSINESS/FUNCTIONAL REQUIREMENTS

Design with the end in mind. Discuss goals with existing analytics users and understand their current challenges and their wish lists. The cloud allows for new capabilities such as near real-time data access, data democratization and next-level analytics with access to detailed data, not just aggregates. Create a plan or a vision statement that highlights being an enabler for all lines of business with the appropriate security controls and tools. Use these goals to align IT and lines of business (LOB) to help set the vision to securely get accurate information to the right people at the right time.

NON-FUNCTIONAL REQUIREMENTS

Often called the “ilities,” take stock of policies related to service level agreements (SLAs) between IT and the business, the security requirements for protecting your data, usability requirements of your end users and many other requirements. Common topics include:

- i. Security
- ii. Reliability
- iii. Performance
- iv. Maintainability
- v. Scalability
- vi. Usability

With the cloud, don't be afraid to include some aspirational requirements. Once you have a solid list, label them as either nice-to-have or must-have. Be wary of paralysis by analysis, and choose an implementation window and development methodology — agile, waterfall, etc. — that is comfortable for your organization. Pay particular attention to your high availability and disaster recovery (HA/DR) requirements. A cloud-built data warehouse can save your organization from having to design elegant but expensive solutions to meet the needs of the business.



Three approaches for implementing your migration

Now it's time to figure out what type of migration would make sense. You have options, which include lift and shift, lift, "improve" and shift, and full redesign. The steps you've taken prior to this stage, such as aligning the migration with your organization, and the technical and business strategy, will drive your chosen migration strategy.

LIFT AND SHIFT

Most would consider this to be the safest and most straightforward way to do a migration. The plan is simple: Everything we do with the existing system should be exactly the same in the new system with minimal changes. A lift-and-shift strategy is a good one if:

1. Requirements are narrowly defined (very few new requirements).
2. Time-to-implementation is critical (you need to get off the old system ASAP).
3. Your new system has all of the features and functions of the old system.
4. Your ecosystem of surrounding tools (ETL, BI, system management) requires minimal or no changes.
5. The migration is not the centerpiece of your cloud migration strategy (see below).

Number five is debatable because many cloud migrations require technical changes and changes to the culture of an organization. If your first initiative doesn't provide more than just the same features and functions as your old system, can your organization view it as a "win"? Sometimes, just showing you can migrate without risk to the business, while improving performance in some way, can be good enough for your stakeholders.

LIFT, "IMPROVE" AND SHIFT

This is by far the most popular approach and can help bridge the gap between the old world of data warehousing and the new era of big data. The concept is simple. As you're converting assets, look for opportunities to streamline or improve the data pipeline, how data is organized, when data is transformed and how data is accessed. Then, find ways to take advantage of new capabilities/functions in the system to which you're migrating. The theory isn't to change any of the core functionality of the system but simply take advantage of the opportunity to simplify or streamline.

The benefit here is to show some improvements over the existing process without breaking things and introducing too much risk. The executive team can





use the migration as a proof point to the business. Even though this represents a major shift in IT philosophy, it does so without negatively impacting performance and provides additional business benefits, such as:

1. Faster access to more data
2. More granular data analytics
3. Better performance on individual queries/reports
4. No contention for near-unlimited computing resources

RE-DESIGN/RE-ARCHITECT/ CONSOLIDATE

Many organizations do not have an enterprise data warehouse or data lake. In some cases, they've been disappointed by their attempt to create one. Their data sits in multiple, on-premise systems: some used for OLTP, some used for OLAP and some data sits in file systems just waiting to be analyzed. Changing platforms is viewed as an ideal time to re-architect, or architect for the first time a fully functional data platform capable of scaling with the business.

The platform must meet the requirements outlined in the above sections: handles multiple types of data, and allows end users to use their favorite tools and language (SQL) to create a data democratization strategy for all lines of business. In some cases, the cloud data

warehouse even creates a new opportunity for modern data sharing across and outside an enterprise. Projects such as these are a great way to consolidate infrastructure and get a better handle on contracts, security and shadow IT, while producing incremental results for the business.

Most consolidation efforts start by combining multiple data sets onto the new cloud platform to show their analytic value. The next step would be to gradually restrict access to the legacy systems, while growing the capability of the cloud data warehouse and establishing quick wins. Many of these initiatives also focus on creating new revenue streams, shrinking or eliminating data pipelines and consolidating disparate data repositories.



PART 2

Planning your migration

Executing these steps, and in this order, is not a necessity. Depending on the scope of your migration, you may need more or fewer steps. The key is to design a framework and core elements of the plan that you can work from. Assess your internal skill sets, don't be afraid to leverage the best practices outlined by strategic vendors, and consider partnering with migration experts.

STEP 1: DETERMINE THE SCOPE

Stating the obvious, no two migrations are the same and rarely is the end state well understood. The goal is to create a plan that aligns with the goals of the business, provides capabilities in the shortest reasonable timeframe and sets you on the path for incremental improvement. Your end state could be getting a single workload into the cloud within one month, or it could be migrating your entire analytics platform by the end of the year. It's reasonable to plan for a one-year ROI, which you can even accelerate under certain scenarios.

STEP 2: DOCUMENT THE "AS IS"

This isn't the most glamorous part of a migration but it's likely one of the most critical. You'll need to communicate both internally and externally, and up and down the reporting chains, regarding the current "as is" implementation. A short list of assets to migrate include but are not limited to:

1. All sources that populate the existing systems
2. All database objects (tables, views, users, etc)
3. All transformations, with schedules for execution, or triggering criteria
4. A diagram from the interaction of systems/tools

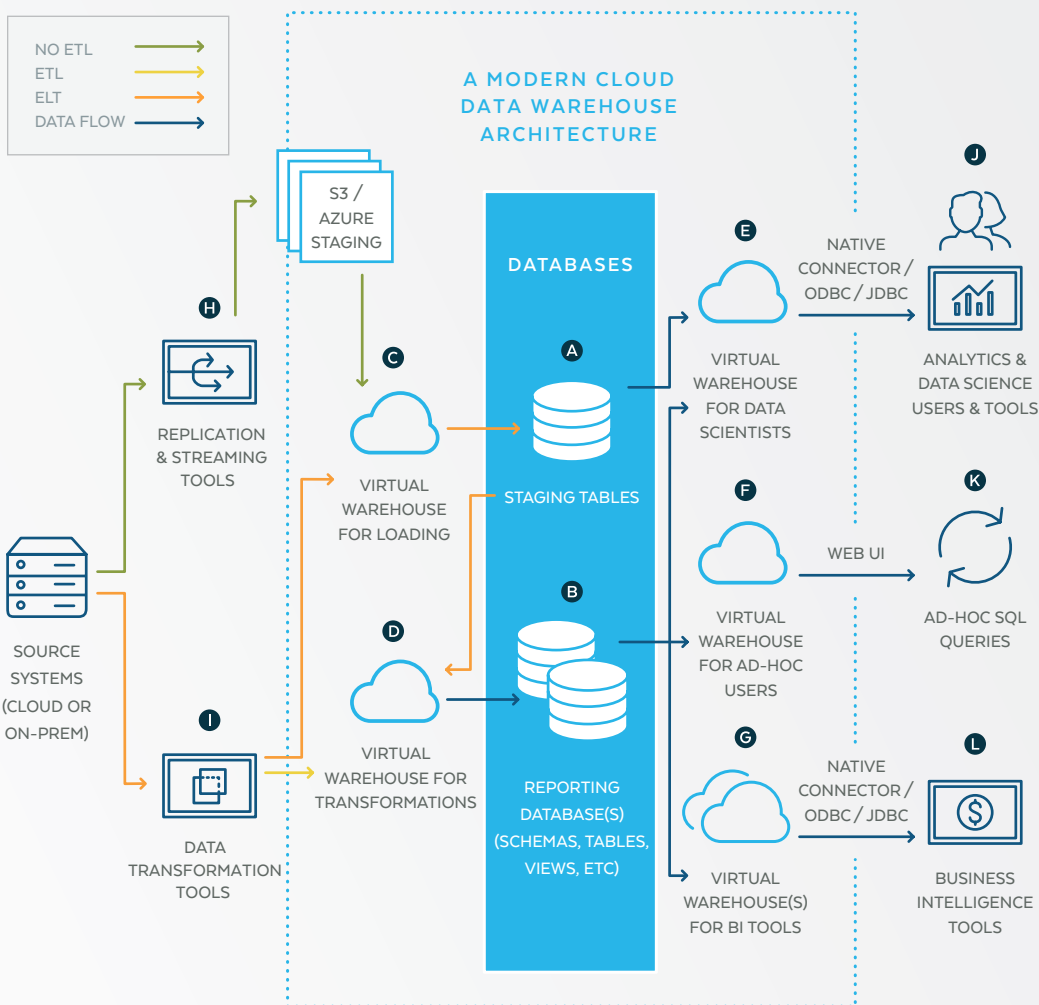
STEP 3: DETERMINE THE APPROACH AND ASSEMBLE THE IMPLEMENTATION TEAM

Multiple options exist here. We've outlined above the most common approaches but there are combinations of these. You could choose to implement one method to get to initial capability and another as you approach full production. Creating high-level milestones at this step is a good way to segment when a capability will be available, and which requirements you'll satisfy via release schedules.

STEP 4: DIAGRAM THE “TO BE”

Once you get your arms around what you want to migrate and when you want capabilities available, you can begin to document the “to be” architecture. Be aware, there is no “one chart to rule them all”. You’ll be communicating to technical,

business and executive audiences, and each will want different levels of detail on the initial operating capability (IOC) and final operating capability. Below is a “to be” example of modern cloud data warehouse architecture:



STEP 5: PLAN FOR YOUR DATA LOAD AND “SIZING”

One of the most challenging aspects of migrating to the cloud is moving data and changing your paradigm to take advantage of the elastic resources of the cloud. There are three pieces to the puzzle:

1. Initial load

The initial load can be challenging based on data volumes and security requirements. Work closely with your security team, and the lines of business that own the data to make sure you don't have to go through a tokenization/obfuscation process before moving data into the cloud. Many organizations segment their data into “zones” or “layers” inside their data lake: raw, curated, aggregated and cleansed areas. Then make the decision which data sets are okay to move, taking into account regulatory and data privacy compliance standards such as PII, PCI and HIPAA. Regarding the volume of data, as networking gets better, this issue starts to go away as you can move terabytes of data into the cloud.

TIP: Many organizations receive data externally from partners and vendors (Salesforce, etc). It might be prudent to dump these data sets into a cloud object storage service to keep it from becoming classified as an on-premise asset. If data is coming over the Internet, it should be ok to secure it in the cloud.

2. Ongoing updates

Each source of data, the ETL logic and integration with the data lake strategy, will dictate the methods used for updating data in your cloud enterprise

data warehouse. Some sources and ETL tools support change data capture (CDC) strategies. Others might support all inserts, while others might require a full refresh. There is no one size fits all approach. If you have the requirement that all data must end up in the lake, you can either write to the lake first or the data warehouse first. It's your choice. This is a place where you could limit the amount of change in the existing process and revisit after the initial implementation.

3. Planning for warehouse usage and storage

Typically, most organizations execute a POC and go through an ROI exercise before executing a migration. At this phase, it's usually a good idea to re-validate the usage plan and work the operational side of the equation with regards to how to monitor availability and how to govern usage of the system. Because some cloud data warehouses provide the ability to scale up and down, turn resources on and off, segment workloads and resources, and auto-scale both processing power and storage, the model changes from time slicing a fixed resource (limiting your business user access) to allocating resources based on business need and value. You no longer have to do a big planning exercise to handle your largest workload and leave the system underutilized for the other 364 days of the year.

TIP: If you are integrating your migration with your data lake strategy, be aware of transfer charges of moving data between regions or cloud providers.



STEP 6: CONVERT ASSETS

This step refers to defining data warehouse/database assets you may need to convert. These include data definition language (DDL), role-based access control (RBAC) and data manipulation language (DML) used in scripts. The good news is that most relational databases leverage the ANSI-SQL standard. Most of the changes will revolve around ensuring DATE and TIMESTAMP formats are converted correctly, and the SQL functions used to access those are checked for compliance. (Not all vendors implement functions the same). Some cloud data warehouses simplify DDL by eliminating the need to partition and index, so your DDL becomes much cleaner (less verbose).

STEP 7: SETUP YOUR “TO BE” ENVIRONMENT AND TEST CONNECTIVITY / SECURITY

It should be no surprise that you'll have to complete your networking, proxy and firewall configurations during your migration strategy. It usually helps to have a chart or two outlining what ports and URLs you will need to access. You will also want to work with your security group to download and install any drivers (ODBC, JDBC, etc) or support software such as a command line interface (CLI), which



most DBA-type developers prefer to use when interacting with a modern cloud data warehouse. You will also want to set up your account parameters such as IP whitelisting and role-based access control before opening the environment up to larger groups.

STEP 8: TEST THE PROCESS END TO END WITH A SUBSET OF DATA

Once you have connectivity worked out for all tool sets, it's best practice to test your process from end to end for both functionality and performance. If you're coming from an existing system, you'll have the advantage knowing what your service level agreements (SLAs) are with each line of business and requirements for each step in the process – load, transform/aggregate, query execution, etc. It's always best practice to select test cases and data sets critical to early success. This is also an opportunity to implement the “improve” part of the migration, if that's the methodology you chose.

STEP 9: MIGRATE THE DATA AND TEST PERFORMANCE

Before going live, you'll want to re-run some or all of your performance tests from step eight to ensure the system is configured for individual query performance (size of the warehouse) and for concurrency (scaling out the warehouse). You'll want a champion from each of the critical lines of business to test accessing the system, making sure their tools work and ensuring they get the performance you expect for them. Validate that you're getting the same calculated results from the old to the new system.

STEP 10: RUN YOUR EXISTING AND NEW SYSTEMS IN PARALLEL

When replacing an enterprise data warehouse or data lake, it's easy to dual load the systems and run them in parallel. From there, gradually move users or groups onto the new system, trying not to disrupt business operations. Target groups that expressed the most challenges or concerns during the planning and requirements phase as they will be the most receptive and probably become the most vocal advocates about success.

STEP 11: PICK YOUR CUTOVER DATE

You can run the systems in parallel for a while but once your end users experience the new performance benefits they will never want to go back. You'll likely want to run the systems in parallel for at least one major reporting cycle—a week, a month or a quarter. Once you pick an official cutover date, continue to dual load the systems in case you run into a problem. Once everyone seems happy, it's time to pop the champagne and retire the old system.



Taking advantage of your new, cloud data warehouse

Now that your migration is complete, the work isn't done. It's time to start taking advantage of the new capabilities. One strategy is to look for ways to improve performance and get data into the hands of end users more quickly, or distributed to more users. Some options include:

DETERMINE LOB AND END-USER NEEDS

You compiled an initial list of requirements from the LOBs earlier in the process. Now it's time to revisit that list and start a dialog with the LOBs, educating them on what's possible with the new system. Some of these new capabilities include but are not limited to:

- a. **More data**—access to detailed records and reaching back years, not just months or weeks
- b. **Different data types**—Structured and semi-structured (JSON, AVRO, XML, PARQUET, ORC)
- c. **Cleaner data**
- d. **Better formatting/modeling**—changing schemas from third normal form (3NF) to star or other data models
- e. **Faster performance**—many organizations use summary tables or materialized views to improve performance

ACCESS TO CROSS-BUSINESS UNIT DATA

Along with the LOB specific data sets, most organization also want access to the latest and greatest information from other business units and sources. Some modern cloud data warehouses provide several methods for controlling data access while still enabling curated data sets to other end users. In some cases, organizations share or receive data sets via FTP with other organizations. Moving data this way is time consuming and expensive, especially as data volumes and frequency increases. Some modern data sharing features associated with modern cloud data warehouses eliminate the need to transfer and transform data, streamlining the process and reducing ETL cost and complexity

OPTIMIZE (RETHINK) YOUR LOAD STRATEGY

Typically, enterprise data warehouses have load windows to execute in batch overnight, making yesterday's data available for analysis the next morning. With today's technology, you can load data and query data without contention, opening the possibility to load data 24/7, and providing data access sooner to end users. You can spin resources up and down instantly, which allows you to load data as fast as you want, and at the same price point no matter the compute resources you spin up.

For example, let's say it takes four hours to load 1TB of data every night. In some cases, you're using one node for this workload. With near linear scalability, if you execute the same workload with a two-node cluster, it would load twice as fast but cost exactly the same. Double the cluster size again and it gets done twice as fast again at exactly the same cost.

ANALYZE SOURCES FOR SIMPLIFICATION/STREAMLINING

In some cases, the tools/systems you were using in the past may not be necessary in the future. Today's modern cloud data warehouses can handle new data types and process data more efficiently, blurring the lines between OLTP, OLAP and the data lake. Large organizations tend to have many tools, all bought for specific functions or capabilities.

Now is good time for an ETL/ELT tool rationalization strategy. Some of these vendors have upgraded their products to work in a cloud environment. Do not try to standardize on a single tool. Instead, make technical recommendations for accessing specific sources, or, if you are implementing a data lake strategy, use different tools for accessing different zones or layers.

You should also analyze the tools in the context of how quickly they can move data from your original system to your new platform. Most organizations are moving from extract, transform and load (ETL) strategies to extract, load and transform (ELT) and stream processing strategies to take advantage of on-demand scaling.

By separating compute resources from storage resources, you have the ability to load data without impacting query performance, making near, real-time data processing at scale a reality. You just need to figure out how much change your organization can absorb at one time.

IDEMPOTENT LOADING

Never heard of it before? This isn't a new term, but has been coming up in lots of conversations related to big data processing. The goal is simple. You want to continually load data into your system, but if something goes wrong along the way you don't want to get confused about what has been processed and what hasn't. Idempotent loading means that it doesn't matter if you load a file or record once or ten times, you end up with the exact same result. This is pretty powerful and can be achieved using a combination of the COPY command and transformation using the MERGE commands found in modern data processing platforms. It's especially good for streaming data with updates or deletes but you can implement it in other situations as well.

ZERO-COPY CLONE, TIME TRAVEL AND UNDROP

Some modern cloud data warehouses have the ability to query "back in time", between the previous 24 hours and 90 days. Many organizations use this capability to transform the ELT pipelines. So, if something goes wrong in step 45 of a 50-step process, you don't have to start over. You use time travel to set yourself up to a known state (step 44 that executed

correctly), fix step 45 and continue processing. Most legacy platforms would force you to reload or recover the original table and start over. If you need a snapshot of your database, or table before doing an update, use a zero-copy clone function to make a copy of either but without duplicating the data. That's a big deal, especially if you're currently doing large ETL jobs to pull terabytes of data into data marts for data science or test/QA teams to work with. And finally, did you accidentally drop a production-level database or table while doing the nightly change? No worries. UNDROP works really well in that situation, and magically all of the data and schema reappears.

Conclusion

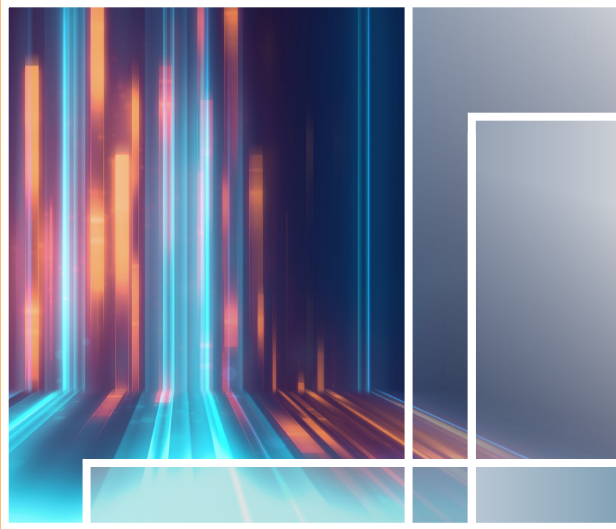
MORE AND MORE DATA IS BORN IN THE CLOUD

Enterprises now realize there are efficiency and performance advantages to storing, analyzing and sharing data in the cloud. This approach removes numerous steps for enterprises, and removes the chaos that can result from siloed, duplicated data that becomes disconnected from its original source. A carefully planned migration can lead to significant advantage over conventional data warehouses, including more capabilities at lower cost.

Learn more

Click [here](#) to get more information about how to modernize your data warehouse and to get instructions specific to migrating from your existing, on-premises data warehouse to Snowflake.





About Snowflake

Snowflake is the only data warehouse built for the cloud, enabling the data-driven enterprise with instant elasticity, secure data sharing and per-second pricing, across multiple clouds. Snowflake combines the power of data warehousing, the flexibility of big data platforms and the elasticity of the cloud at a fraction of the cost of traditional solutions. Snowflake: Your data, no limits. Find out more at [snowflake.net](https://www.snowflake.net).