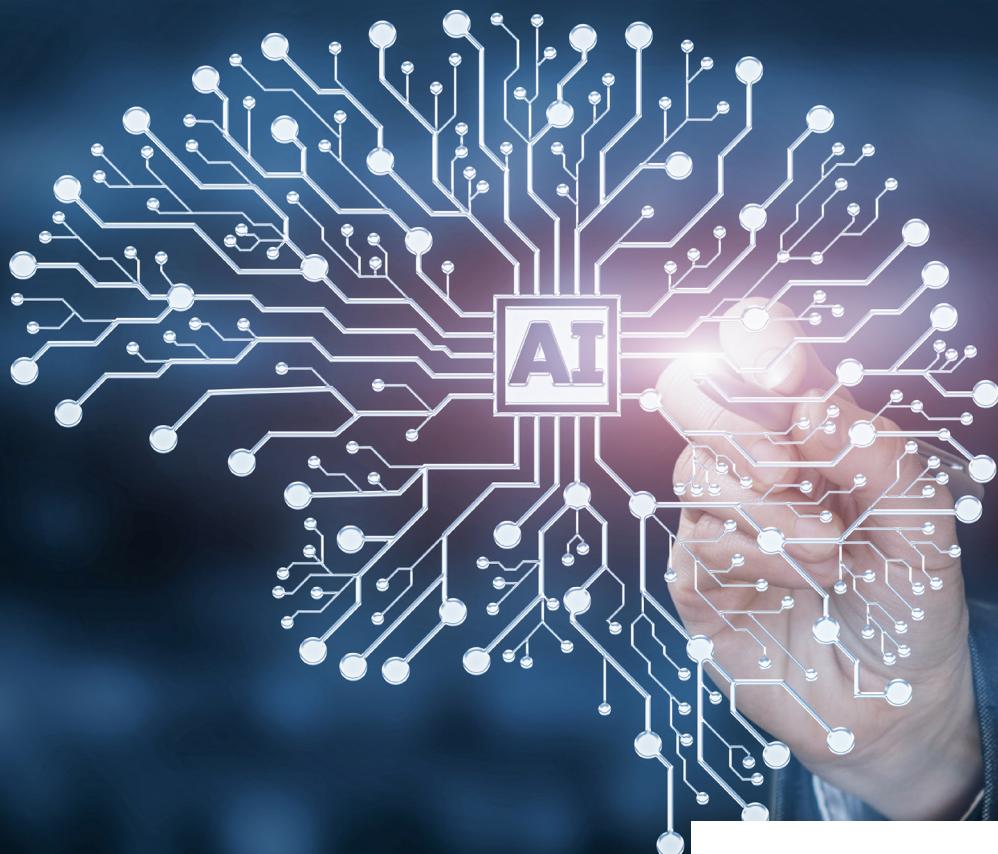# Swiss Banking

# Generative AI in Banking – A Comprehensive Overview

**What do banks need to consider before jumping into the rabbit hole?**

# Executive Summary

Generative AI (GenAI) is one of the most transformative technologies in today's rapidly evolving landscape, offering significant opportunities for the banking industry. By leveraging its capabilities – such as information retrieval, drafting, summarisation, and translation – while understanding its limitations, banks can benefit in multiple ways:

- **Increasing Employee Productivity** – Automating repetitive tasks such as document summarization, report generation, and translation allows employees to focus on higher-value activities.

- **Improving Operational Efficiency** – Automated risk assessments, compliance reporting, and fraud detection streamline back-office processes, reducing cost and errors.

- **Enhancing Customer Experience** – AI-driven chatbots and personalised financial recommendations provide clients with more responsive and tailored interactions.

- **Innovating Financial Products and Services** – AI-generated insights enable banks to create customised investment strategies and credit risk models.

As banks navigate this rapidly evolving landscape, they need to carefully integrate GenAI applications within a highly regulated environment. Given the unique circumstances of each bank – such as size, market position, strategic objectives, and regulatory environment – there is no one-size-fits-all approach to implementing GenAI. In general, successful adoption of GenAI requires banks to act across multiple dimensions: AI initiatives must align with business strategies, be supported by robust IT infrastructure, and adhere to rigorous risk controls. Furthermore, fostering a culture of innovation and providing comprehensive training on GenAI's capabilities and limitations are essential to overcoming adoption challenges.

Against this background, the framework presented in this paper introduces a generic structured approach to guide banks through the integration process, ensuring high-quality outcomes while mitigating risks during execution. It differentiates three task streams – strategic, organisational, and technological – along four phases of GenAI adoption.

- The **exploration phase** focuses on building a foundational understanding of GenAI, identifying potential use cases, and raising awareness among decision-makers. From a strategic perspective, banks must define clear objectives and align GenAI adoption with their broader business goals. Organisationally, fostering a culture of innovation and educating employees on GenAI's capabilities ensures a smooth transition.

- During **analysis and roadmap phase,** banks evaluate the feasibility, prioritise high impact use cases, and create a structured implementation plan. Cross-functional collaboration ensures compliance, while proof-of-concept projects validate technological requirements.

- The **basics and implementation phase** is where the selected GenAI use cases are executed. The Implementation phase focuses on deploying selected GenAI applications by establishing governance, training employees, and integrating AI into workflows. A secure and scalable IT infrastructure ensures regulatory compliance and operational stability.

- Finally, in the **scaling and continuous improvement phase,** banks monitor AI performance, scale successful use cases, and continuously refine processes. Ongoing training, performance optimisation, and cybersecurity enhancements ensure sustainable adoption and compliance.

Looking into the crystal ball, the next phase in AI adoption might most likely be Agentic AI, which enables autonomous decision-making and real-time workflow management. By leveraging structured implementation methods such as prompt chaining and parallelisation, banks will be able to enhance accuracy and efficiency. However, challenges around governance, security, and transparency must be carefully managed.

By addressing strategic, organisational and technological considerations, banks can unlock the full potential of GenAI. This will not only enable them to remain competitive but also position them to meet the evolving expectations of customers and employees while navigating regulatory and legal complexities effectively.

# 1    What's the deal? – Introduction to GenAI

In the rapidly evolving technology landscape, generative AI (GenAI) stands out as one of the most promising and transformative innovations. The advancements of this technology can be attributed to more powerful computers, larger datasets for model training, and ever-enhancing machine learning algorithms.

Various studies have concluded that GenAI will be the next universal technology that could boost productivity growth in many domains of the economy.[1,2]  Compared to other European countries, the Swiss economy could benefit disproportionately from GenAI, as dominant industries in Switzerland with high potential for GenAI applications contribute strongly to GDP. The financial sector in particular, which accounts for roughly 5.4% of jobs and around 9.4% of gross value added in Switzerland, is poised to be one of the biggest beneficiaries of GenAI.[3,4] It is estimated that the introduction of GenAI will impact numerous tasks performed by bank employees. Unsurprisingly, as of 2024, six out of ten bank employees are estimated to already use GenAI in day-to-day business.[5]

At the same time, the technology is subject to numerous technological limitations and requires careful consideration when deployed in a highly regulated industry such as banking. As banks strive to maintain their competitive edge and meet the ever-increasing expectations of their customers, leveraging GenAI while understanding its technical and legal limitations has become essential. Against this backdrop, this paper aims to provide a comprehensive overview of GenAI in the Swiss banking industry by showcasing its general capabilities as well as limitations, demonstrating its relevance for the banking sector along with selected use cases, and pointing out important elements to consider when deploying GenAI solutions in the Swiss banking context.

## 1.1    What is GenAI, and how is it different from traditional AI?

GenAI refers to a class of artificial intelligence models that can generate new content, such as text, images, and audio, based on the patterns "learned" from large amounts of data they have been trained on.[6] Large Language Models (LLMs) are a specific type of GenAI designed primarily to process and generate human language but can also be adapted for tasks involving other formats, such as speech and video, through integration with additional technologies. Traditional machine learning algorithms must be developed and trained on task-specific data to perform one intended task; hence, they are rather narrow regarding the

---

1    🔗 McKinsey, Economic potential of generative AI (2023)

2    🔗 Boston Consulting Group, GenAI Increases Productivity & Expands Capabilities (2024)

3    🔗 strategy&, Schweiz unter Ländern mit weltweit grösstem Wachstumspotenzial durch generative KI, German only (2024)

4    🔗 Swiss Banking, Bedeutungsstudie (2024)

5    🔗 Accenture, The age of AI: Banking's new reality (2024)

6    In the context of artificial intelligence, generative refers to the ability of a system to generate new, original data or new content that has similar characteristics to the training data rather than just analysing or classifying existing data. This differs from discriminative models, which are designed to distinguish between different existing categories of data.

leverage of the model (e.g., classification or prediction). GenAI, in contrast, is "pre-trained" on a vast amount of data – hence removing the tedious training procedure. The breadth and depth of pre-trained tasks are broad but do not come without new risks and shortcomings – which will be discussed later in the paper.[7]

| | Generative AI | Predictive AI |
|---|---|---|
| **How** | Uses existing data to generate new variations | Uses historical data to identify patterns and predict future trends |
| **Output** | Text, images, audio, videos, synthetic data | Probabilities, risk scores, trends |
| **Strengths** | Creativity & highly realistic outputs | Improves decision-making, helps with risk assessment, optimises business strategies |
| **Limitations** | Can produce biased and inaccurate content (hallucinations), requires massive computing power, potential copyright issues | Can be biased and inaccurate if trained on poor-quality data, struggles with unpredictable events |

**Figure 1:** Differences Between Generative AI and Predictive AI · **Source:** own contribution, based on TechTarget[8]

One of the most well-known examples of GenAI is the "Generative Pre-trained Transformer" (GPT) series developed by OpenAI. These models have demonstrated remarkable proficiency in generating human-like text, answering questions, and even performing creative tasks such as writing poetry or composing music.

**"GenAI will most likely not replace machine learning but complement it."**

GenAI's versatility and scalability make it a powerful tool for a wide range of applications across different domains. As such, it will most likely not replace machine learning but complement it. Traditional machine learning will continue to be used for very specialised tasks, but it remains training and data intensive. GenAI, in contrast, can be expected to deliver value for more mundane "day-to-day" tasks without the need for extensive training due to LLMs' four fundamental features:

1. **Versatility in language tasks:** LLMs can perform a wide range of language-related tasks, including text generation, translation, summarisation, sentiment analysis, and more, often without requiring task-specific training. This versatility makes them valuable tools for many applications, from content creation to customer support.

2. **High-quality text generation:** LLMs can produce human-like text that is coherent, contextually relevant, and often indistinguishable from text written by humans. This capability is particularly useful in areas like creative writing, automated report generation, and conversational AI.

---

7        McKinsey, What is generative AI? (2024)

8        TechTarget, Generative AI vs. predictive AI: Understanding the differences (2024)

3. **Efficient transfer learning:** LLMs are pre-trained on vast datasets and can be fine-tuned with relatively small amounts of task-specific data to perform well on new tasks. This efficiency reduces the need for large, domain-specific datasets and extensive training from scratch, saving time and resources.

4. **Scalability and adaptability:** LLMs can be scaled up with more parameters, additional training data, and more refined technologies to improve their performance further. They can also be adapted to various industries, including healthcare, finance, education, and law, where they can assist with specialised tasks such as medical diagnosis, legal document analysis, and personalised learning.
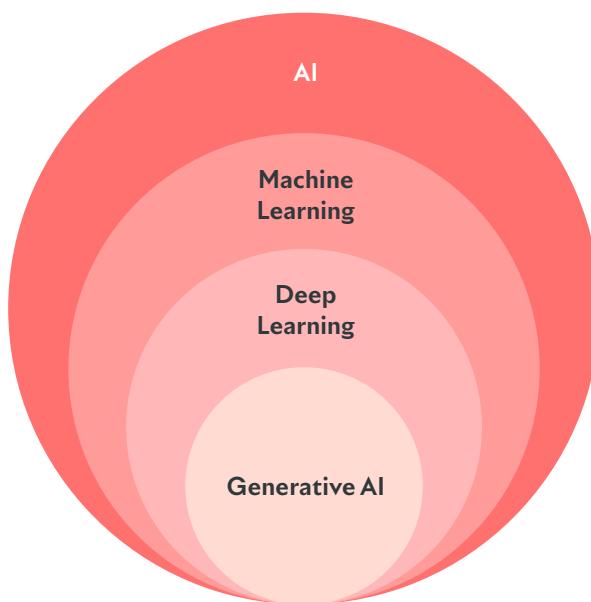


**Artificial Intelligence (AI)**
A discipline (branch) of computer science that deals with systems that can reason, learn and act autonomously.

**Machine Learning (ML)**
ML, a subfield of AI, is a program or system that trains a model on input data and predicts outputs on unseen data. Common classes in ML are supervised, unsupervised, semi-supervised, and reinforcement learning.

**Deep Learning**
Deep learning uses Artificial Neural Networks – allowing them to process more complex patterns than traditional ML.

**Generative AI**
Generative AI is a subset of deep learning. It uses artificial neural networks and can process both labelled and unlabeled data using supervised, unsupervised, and semi-supervised methods.

**Figure 2:** Categorization of AI[9]

GenAI introduces significant operational risks beyond those of predictive AI. While substantial, these risks can be mitigated as outlined below. Note that this list represents only the bank's perspective and is not exhaustive.

- **Misinformation and deepfakes:** GenAI has the capability to fabricate highly realistic text, images, and videos, which can be exploited to create deepfakes or spread misinformation. In banking, this could manifest in various ways: fake emails or video messages from bank executives to employees or from fake clients to Relationship Managers, instructing them to wire money to fraudulent recipients illegitimately; falsified loan documents or fabricated financial reports submitted by credit applicants; or audio files mimicking clients to circumvent voice authentication in online banking environments. Such scenarios could result in security breaches, as well as financial and reputational damage.

---

9    &#x1F517; Based on: Alzubaidi, L., Zhang, J., Humaidi, A.J. et al. Review of deep learning: concepts, CNN architectures, challenges, applications, future directions (2021).

- **Hallucination:** GenAI is not bulletproof. As it's probabilistic in nature, it is known to create responses that seem credible but are false. This issue, known as "hallucination", presents a serious risk for maintaining trust and accuracy in financial services. Moreover, many LLMs do not specifically mention sources, creating difficulties in verifying the accuracy of responses. And even if asked to provide sources to back up responses, some of them may be made up.

- **Intellectual property and copyright infringement:** GenAI's ability to produce content that mimics existing works raises significant concerns about intellectual property (IP) and copyright infringement. In banking, this could, for example, involve the generation of research reports or financial models that closely resemble proprietary tools or data from other institutions. Determining the ownership of AI-generated content becomes more challenging, leading to potential legal disputes over financial products or reports. Traditional AI, which typically works within predefined datasets and models, poses fewer risks in this area, making the IP challenges with GenAI unique to its content-creating nature.

- **Bias and ethical concerns:** GenAI models, by generating new content, can unintentionally perpetuate or amplify biases embedded in their training data. This could result in biased credit scoring algorithms, unfair loan approval processes, or biased investment recommendations in the banking sector. For example, a generative AI model might create new customer profiles or financial products that unintentionally reflect discriminatory practices, leading to ethical and regulatory challenges. While traditional AI also addresses bias, it primarily interprets existing data patterns rather than generating new biased content. However, traditional AI can also lead to biased human decisions if those decisions are based on insights drawn from biased data. In GenAI, this capability to create and perpetuate bias presents a more complex ethical dilemma, especially in the context of financial inclusion and fairness.

- **Cybersecurity & information security:** GenAI introduces unique cyber security challenges that differ from those associated with traditional AI systems. One such risk is prompt injection, where malicious actors exploit how these models process input prompts. [10] Since they respond to free-form text rather than structured inputs, adversaries can craft deceptive prompts to manipulate their behaviour or gain unauthorised access to sensitive information. Additionally, vulnerabilities like jailbreaks, information leakage, model theft, data poisoning, and supply chain vulnerabilities further underscore the need for a comprehensive security strategy. In essence, the dynamic nature of GenAI systems requires robust security measures to prevent and mitigate these threats.

- **Third-party risks:** Integrating third-party GenAI solutions introduces risks such as unintended exposure, where banks may lack full visibility into how external AI models operate, data leakage, which can lead to the accidental disclosure of sensitive information, and supply chain vulnerabilities, where reliance on external providers creates security gaps.[11]

---

10    See for example 🔗 Vischer, Teil 6: Die andere Seite der Medaille: Wo wir KI vor Angreifern schützen müssen (2024)

11    🔗 Based on: McKinsey, Implementing generative AI with speed and safety (2024)

## Deep Dive: How do LLMs work in detail?

All in all, it's essential to understand that LLMs are still "only" probabilistic algorithms which, at their core, adopt statistics. Whilst they have billions of parameters, they can only perform one task: given a sequence of words, predict the most likely next word. During extensive training, which can last weeks or even months, these models learn a probability distribution over words. While predicting a single next word might seem trivial to humans, it involves complex calculations utilising all the model's parameters. For generating longer sequences, LLMs operate in an autoregressive manner, consuming their own outputs and adding one word at a time.[12]
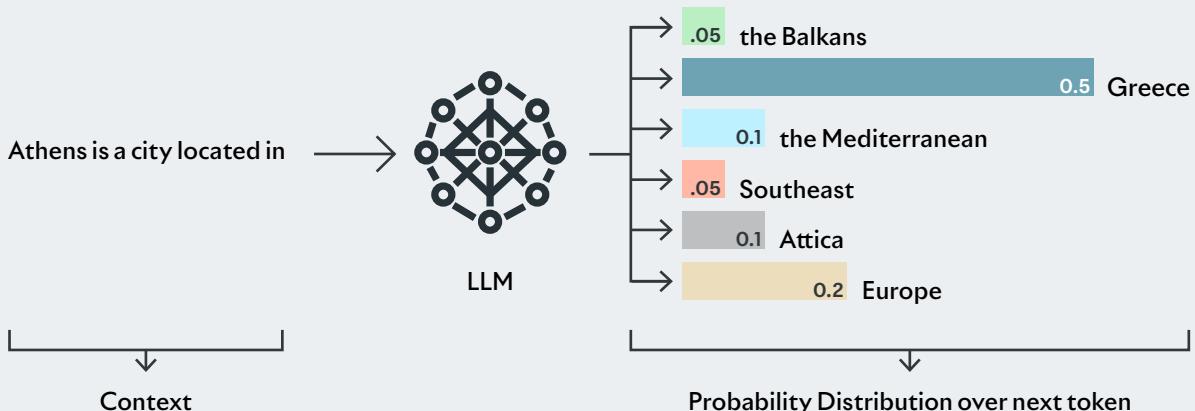


**Figure 3:** Next word prediction of an LLM based on context and the predicted probability of the next token · **Source:** Own contribution based on NVIDIA.[13]

Although these models lack an inner monologue, they can reflect upon their "thoughts" by generating words. Since LLMs can only add new words and not remove them, any incorrectly chosen word (for instance, "yes" instead of "no") will be reinforced by subsequent words – recall that every word is chosen such that it fits best to what precedes it. This may result in so-called hallucinations, or in other words, statements that are persuasive yet fundamentally inaccurate. To address this, it is advisable to ask the LLM for a justification before committing to a final decision, allowing the model to reconsider the context and avoid potential errors. For this reason, many LLMs habitually repeat questions in their own words before answering.

---

12    Technically, LLMs work with so-called "tokens" sub-words or sometimes single characters. Due to language variations such as declination, conjugation, and inflexion, the number of unique words can be large (e.g., [to] "run", [he/she/it] "runs", [... is] "running"). By breaking down words into reusable chunks ("run", "-s", "-ning"), it is possible to represent language with a smaller vocabulary while maintaining the semantic meaning of the root word ("run" is consistent across "runs" and "running"). Identifying effective tokenization patterns is an optimisation problem. Most common LLMs use patterns that are tailored to the English language.

13    NVIDIA, How to Get Better Outputs from Your Large Language Model (2023)

Until the early 2010s, researchers have built language models by manually implementing grammar rules and carefully annotating datasets, with perhaps the most known application of these so-called rule-based language models being grammar checkers in word processing apps. Today, however, training LLMs is largely hands-off, relying on next-word prediction with vast and diverse textual datasets. As a by-product of this objective, these models implicitly learn grammar, syntax, and synonyms while also adapting to multiple languages and acquiring vast amounts of knowledge. Consequently, high data quality is crucial for the adequate training and performance of an LLM.

## 1.2    Why is GenAI important for banks?

As the financial industry undergoes rapid digital transformation, banks are presented with opportunities to innovate and strengthen their competitive edge. While using "off the shelf" LLMs might not be a competitive differentiator, the correct implementation and "grounding" within the banks' specific context surely can be. For many banks, the added value of GenAI is becoming increasingly apparent. While there may not be an immediate need to implement GenAI across all operations, it is crucial to incorporate it into mid-to-long-term strategic planning. GenAI offers transformative potential across several key areas, making its integration into a bank's strategy a forward-looking decision – always with the understanding that GenAI enhances human decision-making, creativity, and customer service.

One of the strategic benefits of integrating GenAI into a bank's operations is the general technological maturity of the organisation it fosters, particularly by empowering employees to evolve with the technology. By starting with the integration of GenAI now, banks can gradually build internal expertise and provide a climate and culture in which experimentation with this new technology is encouraged.

**"Integrating GenAI into banking is not about rushing to adopt the latest technology but about thoughtfully incorporating it into mid-to-long-term business strategies."**

GenAI may also aid in tasks that humans are not ideally suited for. The amount of data generated from daily banking operations is growing exponentially. Whether it is transaction data, customer interactions, or market analytics, the sheer volume and complexity of this data are becoming increasingly challenging to manage. Machine learning and generative AI provide the necessary support to process, analyse, and extract valuable insights from vast datasets at speeds and scales that are beyond human capability. Furthermore, it enables the use of unstructured data in new ways. By gradually incorporating GenAI-based tools into data management strategies, banks can ensure that their employees are equipped with the most accurate and relevant information, enabling them to make better, faster decisions.

In a rapidly evolving, data-driven world, integrating GenAI into banking is not about rushing to adopt the latest technology but about thoughtfully incorporating it into mid-to-long-term business strategies.

By leveraging this technology, banks can achieve higher productivity levels, offer superior customer experiences, and stay ahead in a competitive market. However, the successful implementation of GenAI requires an approach that addresses various challenges such as data privacy, regulatory compliance, risk and ethical considerations. By understanding the full scope of GenAI's capabilities as well as limitations and preparing to navigate its challenges, banks can unlock the true potential of this technology.

# 2    How to benefit? – Generic use cases in banking

Enhanced competitiveness, improved client acquisition and retention, increased operational efficiency, greater technological maturity, and the ability to effectively manage vast amounts of data are some of the many benefits that make a compelling case for a careful and strategic integration of GenAI. In this report, we identified several use cases and examples along four generic categories to showcase the benefits of GenAI for banks and their employees.

## 2.1    Increasing employee productivity

Off-the-shelf tools like ChatGPT or Claude can improve an employee's productivity and effectiveness. Some examples of GenAI applications in this context include searching or summarising documents, reviewing and improving the quality of texts, translation, research tasks, and writing code. These tools address a wide area of use cases and everyday tasks. Enabling employees to do more in less time, empowering them to leverage the new way of work, and overcoming digital friction should be one of the top priorities of banks. Use cases like meeting transcription, creation of documents and presentations or searching along personal data spaces (emails, files, etc.) enable productivity boosts and help employees focus on their core tasks. GenAI can help with task prioritisation and workflow automation, allowing employees to focus on what is essential for a bank and reduce time spent on mundane tasks. In addition, accessing internal unstructured data using a company-specific Retrieval Augmented Generation (RAG) tool helps to find and use internal knowledge in a shorter time. Furthermore, GenAI could automatically generate and summarise financial reports, contracts, and compliance documentation, reducing manual effort and minimising errors. In a self-reinforcing loop, increased productivity might also help to improve customer experience: as employees deliver personalised, efficient, and error-free services, it positively affects customer satisfaction.

## Use Case: Knowledge management at the Pictet Group

Beyond productivity improvements that LLMs deliver out of the box, one important application of LLMs is the ability to search through large custom knowledge bases proprietary of a bank. For instance, Pictet Group has deployed an internal chatbot that can answer many questions about the Group's history, HR processes, people directory, group directives and policies, or even IT-related questions. Using the RAG method, the chatbot retrieves the most relevant elements from the knowledge base, which is connected to Pictet's document libraries, to answer specific questions and provide sources and links to the relevant documentation. In addition, the chatbot can offer contextualised answers based on the employee's access rights to specific resources, i.e., enhancing answers by drawing on information that is only visible to the employee.

## Use Case: Generative AI Empowerment Program at Raiffeisen

A study at Raiffeisen Bank revealed that the successful implementation of generative AI applications depends on employees' ability to craft effective prompts. Prior to deploying sophisticated solutions like Microsoft Copilot, the financial institution launched Copilot Chat in March 2024 – a more lightweight version of Microsoft Copilot. This approach aims to familiarise personnel with AI interaction in a chat environment, develop their prompting expertise, and prepare the workforce for more advanced tools in the future. This strategy also provides additional time to develop integrated platforms for internal retrieval-augmented generation (RAG) systems.

Although Copilot Chat is regularly used, it has not yet become seamlessly integrated into daily operations. To expand adoption across the bank, they implemented two initiatives: First, they appointed approximately 300 staff as local AI champions who facilitate communication between departments, regional branches, and the AI project leadership. Second, they established multiple learning channels to enhance prompting capabilities, including beginner-focused 10-day training sessions, a collaborative Teams space with educational resources, weekly prompt challenges and missions, and company-wide "Promptathons" – engaging events where participants discover innovative applications of AI technology. This comprehensive enablement strategy yielded impressive results, driving a 50% increase in Copilot Chat adoption. More staff members now leverage AI technology to enhance their productivity. Given this success, the institution is optimistic that its workforce has developed sufficient expertise to handle more sophisticated AI applications in the future.

## 2.2    Improving operational efficiency

GenAI can be applied to automate entire parts of processes and lead to considerable efficiency gains. Automating data entry and validation tasks can streamline document generation, compliance reporting, and customer onboarding. Also, it can review regulatory requirements, monitor compliance, and support real-time risk assessment, helping employees focus on strategic tasks over manual checks. In the context of legal and compliance and other essential support processes, work often involves structuring and analysing large amounts of information. GenAI systems are efficient helpers for such time-consuming tasks. In loan processing, GenAI can support risk assessments and documentation, speeding up decision-making and improving case management.

### Use Case: Using corporate language for translation at Julius Bär

At Julius Bär, GenAI supports the translation of corporate content, ensuring communication aligns with the bank's tone and terminology. A fine-tuned language model helps meet these needs but requires careful calibration to avoid underfitting and overfitting. Underfitting arises when the training dataset is too small or lacks diversity, leading to generic translations that miss nuances. For instance, the model might fail to capture specialised financial terminology or the required tone for client materials. Overfitting, however, occurs when the dataset is overly specific, causing the model to struggle with new or varied content. To address this, the team trains the model on human-produced content that adheres to the desired standard and sets the ground-truth which the LLM should thrive for. Quantitative as well as qualitative testing and feedback collection ensure that the model adapts to evolving translation needs. This, and always keeping the human in the loop, allows Julius Bär to use GenAI effectively while maintaining high standards for clarity and accuracy.

## 2.3    Improving customer experience

Banks are increasingly deploying GenAI to improve customer experience by leveraging its key features. These include round-the-clock availability for customer support with automated responses to customer queries, as well as the ability to deliver personalised offers and communication.

By analysing customer data and preferences, AI-powered chatbots and virtual assistants can aid in delivering tailored financial advice, answer questions, and assist with transactions, increasing customer satisfaction and loyalty. Moreover, it also aids in creating targeted marketing campaigns to enhance customer loyalty while analysing customer data to identify opportunities for cross-selling and up-selling. In direct interaction with customers, GenAI can be used for process-supporting tasks that have no influence on processes with a financial focus. However, GenAI is far from ready for autopilot-like working methods. Therefore, these results must always be subject to human scrutiny if they evoke financial decisions.

**Use Case: Enhancing service quality at SIX**

To improve service quality at SIX, generative AI is used to transcribe calls, analyse content, and identify problematic interactions. This approach supports SIX in providing agents with targeted training recommendations and identifying potential business opportunities. The solution is implemented entirely on-premises, utilizing a combination of different open-source models contextualized with proprietary knowledge to extract insights from customer interactions, thereby promoting operational efficiency.

## 2.4    Enhancing existing products and services

GenAI enables banks to develop new financial products and services that meet evolving customer needs. For instance, AI-supported investment advisory can be used to provide personalised investment recommendations based on real-time market data and individual risk profiles. Additionally, GenAI can further enhance the quality of customer advisory through AI-enabled feedback and usage analysis. While GenAI-based tools can significantly improve the quality of banks' client advisory, it is recommended not to allow those tools to interact with customers unchecked - at least at the current stage of development.

However, GenAI can act as an essential preparatory worker for product development in the products and services business area. For instance, so-called AI code completion tools have been central in boosting software developers' productivity long before the meteoric rise of ChatGPT. New, more sophisticated GenAI platforms can not only make seasoned professionals bring new financial products faster to market, their ease of use can also enable employees without programming skills to leverage the power of coding.

# 3    How to succeed? – Considerations for creating optimal GenAI conditions in your bank

To successfully implement GenAI tools, banks must act across multiple dimensions, ensuring that AI initiatives are aligned with business strategy, governed effectively, supported by the right mindset and culture, underpinned by robust IT infrastructure, and secured by rigorous risk controls. Addressing these areas comprehensively will allow banks to harness the power of GenAI while minimising risks and ensuring compliance with regulatory and ethical standards. At the same time, there is no single right path to realising GenAI's potential as a financial institution. After all, some organisations might have already carried out some fundamental tasks while others haven't. As such, the following framework serves as a point of reference for banks who want to implement GenAI to ensure sufficient quality and mitigate potential risks during the project execution. The framework differentiates three different task streams along four phases.

> **"The following framework serves as a point of reference for banks who want to implement GenAI to ensure sufficient quality and mitigate potential risks during the project execution."**

**Three generic task streams**

- **Strategy** considers all tasks that revolve around questions of strategic alignment, use case identification and prioritisation, and the establishment of a long-term vision on how to benefit from this technology.

- **Organisation** considers all tasks that revolve around aspects of people, processes, or policy. They are the glue that keeps strategy and technology aligned.

- **Technology** considers all tasks that revolve around everything related to information technology and data infrastructure, down to the bits and bytes which make up a GenAI application.

**Four generic phases**

- First, there's the **exploration phase,** where banks build a foundational understanding of GenAI and learn to identify opportunities where GenAI might offer benefits. It is crucial that the bank understands the potential of generative AI and strategically sets ambitions in this regard. Raising awareness among decision-makers to address the possibilities is essential. Typically, it is also advisable to tackle internal use cases through concrete exploration.

- Second, there's the **analysis and roadmap phase.** Here, banks assess the feasibility of GenAI applications, prioritise use cases, and create structured project management plans. Proof of concepts should be conducted with broad involvement from various departments to initiate learning within these areas and to adhere to legal and regulatory requirements. This includes IT development, IT operations, architecture, compliance, legal, data protection, risk, and departments with initial application use cases.

- Third, banks tackle the **basics and implementation** of previously identified use cases. In this phase, participants of a GenAI project establish the underlying infrastructure and its governance and then execute the broad implementation of a generative AI tool.

- Fourth, after implementing the GenAI use case, there's the **scaling and continuous improvement phase.** After closely monitoring the performance of the newly implemented application, it's necessary to scale the use case and put effort into maintaining the quality of the GenAI implementation. During the scaling of productive use cases, the bank learns the most.
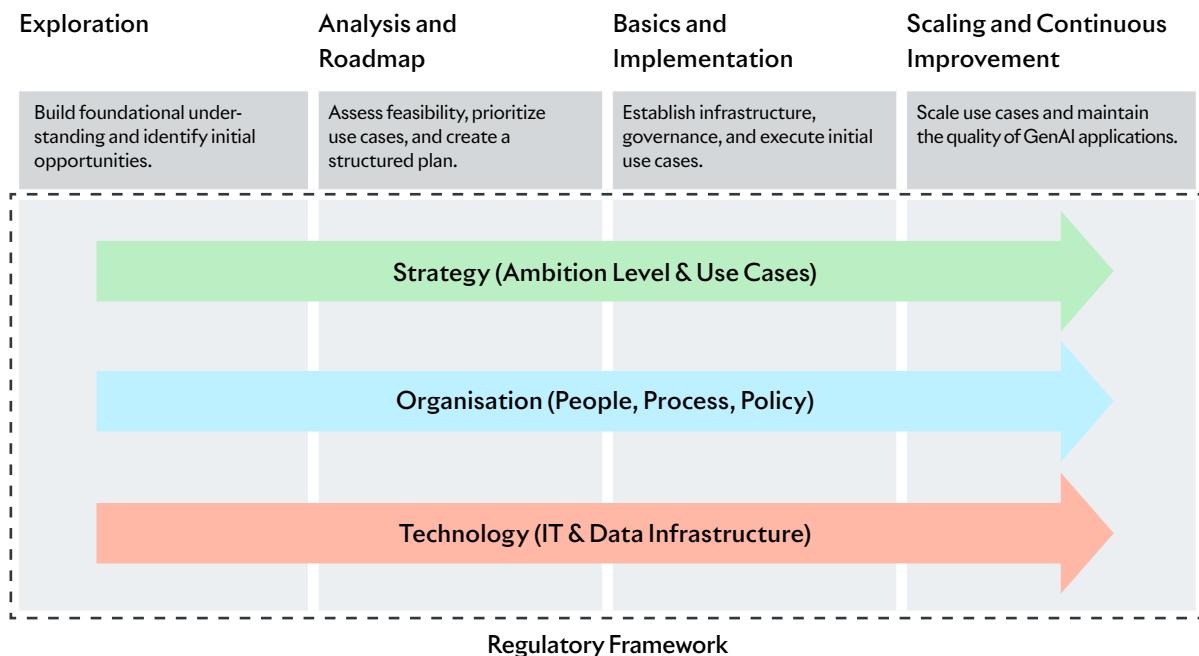
**Framework for Implementing GenAI Projects in Banks**

| Exploration | Analysis and Roadmap | Basics and Implementation | Scaling and Continuous Improvement |
|---|---|---|---|
| Build foundational understanding and identify initial opportunities. | Assess feasibility, prioritize use cases, and create a structured plan. | Establish infrastructure, governance, and execute initial use cases. | Scale use cases and maintain the quality of GenAI applications. |

Strategy (Ambition Level & Use Cases)

Organisation (People, Process, Policy)

Technology (IT & Data Infrastructure)

Regulatory Framework

**Figure 4:** Generic Framework for Implementing GenAI Projects in Banks · **Source:** own contribution

| Exploration | Analysis and Roadmap | Basics and Implementation | Scaling and Continuous Improvement |
|---|---|---|---|

Strategy Alignment

Use Case Identification and Prioritisation

Roadmap Development

Long-term Vision

Strategic Re-Alignment

Organisational Culture

Awareness & Education

Proof of Concept

Impact Assessment

Governance & Compliance

Risk Management

Iterative Scaling

Resource Planning

Pilot or MVP

Infrastructure and Tools
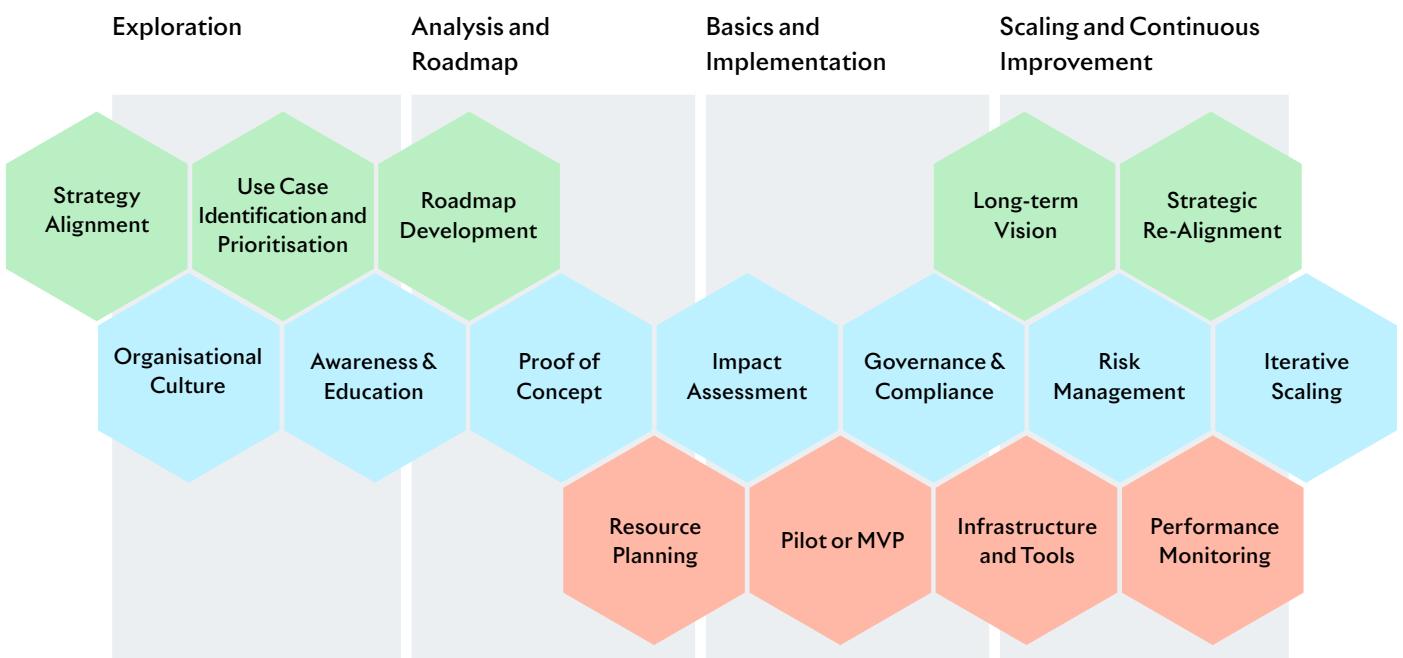
Performance Monitoring

**Figure 5:** Detailed Building Blocks along the Framework for Implementing GenAI Projects at a Bank · **Source:** own contribution

## 3.1    Strategic considerations



To effectively implement GenAI, banks ideally establish a clear strategy that aligns with their business objectives, customer needs, and long-term goals. A clear strategy and long-term vision ensure that the implementation of GenAI tools fits the bank's existing structures and provides guidance whenever an organisational or technical decision needs to be made. Strategic considerations should include the following five tasks:

- **Strategic Alignment:** The AI strategy should be integrated with the bank's overall business strategy, ensuring that (Gen)AI initiatives support core objectives like increasing employee productivity, improving operational efficiency, enhancing customer experience, or expanding product and service offerings. It is paramount that project management define the ambition levels and secure top management commitment. At this stage, the bank sets long-term strategic goals for generative AI integration and determines how the achievement of these goals is measured. Also, the bank should ensure that AI is not seen as the sole responsibility of the IT or data science teams. Collaboration between departments— such as marketing, compliance, risk, and customer service—will ensure that GenAI initiatives are aligned with diverse business needs.

- **Use Case Identification and Prioritisation:** This task involves brainstorming sessions with diverse teams to identify potential use cases (e.g., customer service automation, internal knowledge retrieval). Use cases can be assessed based on dimensions like value creation, cost reduction, revenue increase, and risk mitigation. Furthermore, use cases should be prioritised based on their feasibility, which depends on data quality, technical infrastructure, and employee readiness.

- **Roadmap Development:** Banks should draft a phased roadmap outlining short-term, mid-term, and long-term performance goals. These performance targets for GenAI applications should be measurable and integrated into the AI strategy as well as into broader business objectives.

- **Long-Term Vision:** Once the GenAI application has been established, management should discuss their long-term vision of GenAI use. With the newly gained experience tackling GenAI projects, the bank can explore more advanced applications, such as personalised financial advice and fraud detection.

- **Strategic Re-Alignment:** After having made the first experiences with the new GenAI application and having defined the long-term vision, it is crucial to re-align the GenAI implementation with the overarching business strategy. The initial strategy was fuelled by a forward-looking vision, but as the bank gains hands-on experience, strategic priorities may shift. This phase ensures that lessons learned from real-world deployment are integrated into the AI roadmap, allowing for necessary course corrections. Strategic re-alignment involves evaluating whether GenAI initiatives still deliver value, remain technically and operationally feasible, and continue supporting key business objectives. It may also include refining governance frameworks, updating performance benchmarks, and identifying new strategic opportunities.

## Best Practice: How to identify and prioritise use cases?

To identify and prioritise meaningful use cases for GenAI in a bank, a structured approach from ideation over prioritisation to implementation is recommended. The processing should follow the so-called 3-2-1 approach: Three (3) ideas should be condensed into two (2) feasible use cases, from which one (1) implementation project should be prioritised. It is crucial that the evaluation addresses different dimensions and does not solely focus on a technical perspective. Ideally, the evaluation is conducted by a committee involving various departments to ensure broad anchoring. To ensure a smooth transition and build expertise, it is advisable to implement GenAI-related initiatives within domains where the core business remains unaffected. This approach allows the bank to experiment, learn, and refine processes without jeopardising critical operations. For initiatives involving core banking systems, the introduction of GenAI should initially run passively, with human oversight at every stage. This "human-in-the-loop" model ensures that risks are mitigated while gaining valuable insights. Starting small, with limited scope, is recommended before scaling up to tackle larger, more complex projects.
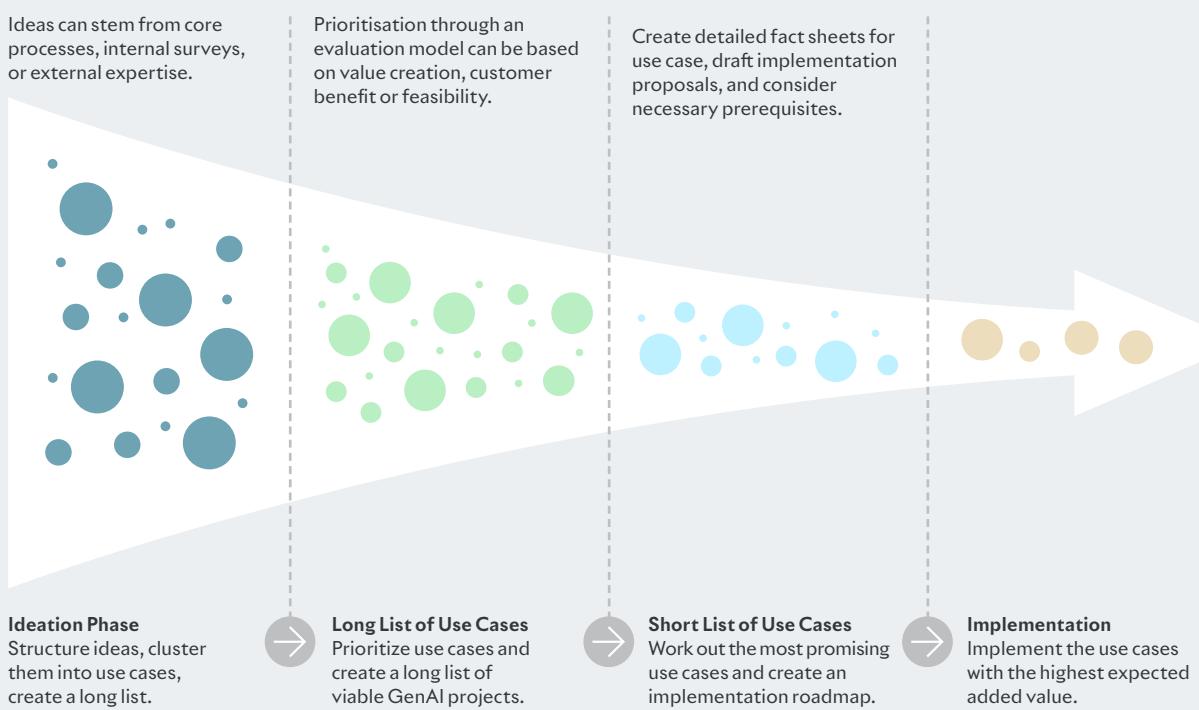
Ideas can stem from core processes, internal surveys, or external expertise.

Prioritisation through an evaluation model can be based on value creation, customer benefit or feasibility.

Create detailed fact sheets for use case, draft implementation proposals, and consider necessary prerequisites.



**Ideation Phase**
Structure ideas, cluster them into use cases, create a long list.

**Long List of Use Cases**
Prioritize use cases and create a long list of viable GenAI projects.

**Short List of Use Cases**
Work out the most promising use cases and create an implementation roadmap.

**Implementation**
Implement the use cases with the highest expected added value.

**Figure 6:** General approach to identify and prioritise GenAI use cases · **Source:** Raiffeisen

The following two dimensions can serve as a point of reference for setting up an evaluation model to prioritise the identified use cases. The model should be adapted to the respective strategy, ambition, and capabilities of the individual bank.

**Dimension 1: Value creation**

- **Scalability:** Number of internal and external users who could benefit from deployment.

- **Yield Increase:** Potential to conclude further business through the particular use case

- **Avoidance of Risks:** Supporting the user by averting legal, regulatory and reputational risks

- **Cost Reduction:** Potential for cost reduction through risk mitigation or efficiency savings

- **Qualitative Added Value:** Significant benefits for the bank and customers through innovation, better user experience (UX), reputation (lighthouse effect), and competitiveness

**Dimension 2: Feasibility**

- **Data Availability:** Assessment of available data quality and effort in provision.

- **Technical Feasibility:** Implementation regarding architecture, existing technology, and GenAI patterns.

- **Internal Maturity:** Employee acceptance, process adaptation, availability of know-how

- **Regulation/Customer Acceptance:** Legal and regulatory restrictions and expected customer feedback

- **Content Reliability:** Use case robustness, explainability, security, fairness, and scope of non-sensitive topics

The gathered information can be summarised in detailed factsheets for each use case. These can serve as a basis for further internal alignment and implementation along a defined roadmap. They can contain information on the following aspects:

1) Summary of the GenAI use case (problem definition, solution approach)

2) Added value, effort, and delivery objects

3) Risks

4) Investment calculation

5) Performance measurement

6) Other key data (e.g., key stakeholders, implementation period, rating from evaluation)

## 3.2   Organisational considerations



Organisational considerations include the three aspects around people, processes, and policies and span all four phases of a GenAI implementation project:

- **Organisational Culture:** A data-driven culture will accelerate the adoption of GenAI by making it an integral part of business processes. If employees are not on board with implementing GenAI applications, the tool will most likely miss its target. A bank should provide continuous training to upskill employees not only on using AI tools but also on identifying use cases where GenAI can offer value. Finally, banks should also consider their talent acquisition and retention strategies. As banks increasingly integrate GenAI tools and applications into their operations, the demand for AI, data science, and machine learning professionals will surge. A clear strategy can help to identify the necessary capabilities in the bank and hire new or re- and upskill existing workforce to align with these.

- **Awareness and Education:** A bank can host workshops and training sessions to introduce GenAI to employees. They should be educated on AI capabilities, associated risks, and their potential impact on banking operations. Moreover, employees must not only understand but also embrace these advancements to ensure successful adoption.

- **Proof of Concept (PoC):** Once employees are up to speed on GenAI risks and opportunities and a GenAI implementation has been put on the agenda, a team can test initial use cases such as off-the-shelf tools (e.g. chatbots for FAQ or document summarisation) or RAG techniques for internal unstructured data exploration.

- **Impact Assessment:** In this task, a bank should analyse how generative AI will affect existing processes, employees, and customer interactions. Experts should evaluate regulatory and compliance implications (e.g., data protection, banking secrecy).

- **Governance and Compliance:** A robust governance framework is essential to ensure that GenAI initiatives are deployed securely and in compliance with existing legal and regulatory requirements.

- **Risk management:** Banks should establish a rigorous risk management and governance framework to mitigate the potential risks associated with GenAI deployment. This includes GenAI specific risks like biases, hallucinations, and IP concerns. For this, banks can develop processes for managing AI-related incidents and regulatory audits and can also develop internal policies and guidelines for AI usage (e.g., transparency, and employee training).

- **Iterative Scaling:** Once the PoC has been made and the initial use case has shown promising results, it is time to expand the use case to other departments. It is essential to regularly review and refine AI models to ensure accuracy and compliance.

## Details on governance and compliance

A robust governance framework is essential to ensure that GenAI initiatives are deployed securely and in compliance with existing legal and regulatory requirements. Key elements of governance include:

- **Establishing a robust Data Governance:** Effective data governance is critical when deploying GenAI systems, particularly those trained on corporate or proprietary data. For example, while training a model on existing libraries to enable features like document retrieval or inventory tracking, improper governance could expose sensitive information. If models are trained on data such as payroll spreadsheets, employees might inadvertently uncover confidential details, such as other person's salaries, through clever prompting. Consequently, banks must implement robust data governance frameworks to mitigate such risks and maintain strict identity and access management (IAM) protocols. This trade-off between training models with comprehensive data and safeguarding sensitive information underscores the need for vigilant oversight.

- **Ensuring Clear Accountability:** Define roles and responsibilities for AI governance, ensuring that senior leadership, compliance, legal, data protection, risk, data scientists, and IT teams are aligned in managing GenAI development and deployment. A chief AI officer or similar role could oversee the entire AI lifecycle. Instead of a centralised approach, one could also envision a decentralised approach to accountability in which a Center of Expertise (CoE) governs the bank-wide endeavours, but the pre-defined employees in the divisions and business verticals remain accountable.

> **"A robust governance framework is essential to ensure that GenAI initiatives are deployed securely and in compliance with existing legal and regulatory requirements."**

- **Defining Guidelines and AI Principles:** Develop guidelines that govern the use of AI, focusing on fairness, transparency, robustness and accountability. These principles should be embedded in every phase of AI development, from data collection to model deployment (see also info box on the FINMA risk monitor and the supervisory expectations towards the use of AI).

- **Monitoring:** Implement a governance process for continuously monitoring AI models to assess performance, prevent model drift, and ensure outcomes align with business goals. Establish clear audit trails for accountability, especially in critical functions like credit risk management or fraud detection.

## Deep Dive: Relevant legal and regulatory aspects to be considered when implementing GenAI projects

Swiss law is generally technology-neutral and principles-based. The following overview summarises the most important laws that – typically – need to be considered when implementing GenAI projects in Swiss banks. Depending on a bank's international activities, possible implications of additional applicable regulations must be considered regarding the deployment and scaling of GenAI projects across the organisation.

- **Data Protection:** The Federal Act on Data Protection (FADP) applies if personal data is being processed. In the context of the development, training and use of GenAI models, the processing principles[14] must be followed (e.g., purpose limitation, accuracy) and data security[15] must be ensured, or a justification must be required.[16] Whether or not AI models contain personal data is the subject of a heated debate.[17] This is particularly important regarding the rights of data subjects, such as the right to rectification or the right to forget. Processing which results in a high risk to the data subject's personality or fundamental rights requires a data protection impact assessment (DPIA).[18] In its Factsheet on DPIAs, the Federal Data Protection and Information Commissioner (FDPIC) specifically mentions artificial intelligence as an example of a new technology constituting a risk factor for the assessment of a high risk.[19] Moreover, if a decision is based exclusively on automated processing and has a legal consequence or a considerable adverse effect, the data subject may generally request to express their point of view and to have the decision reviewed by a natural person.[20] In this context, requirements regarding information obligations, data subject rights, data processing agreements, and the disclosure of personal data abroad need to be followed closely. Mainly, FDPIC requires manufacturers, providers, and users of AI systems to make the purpose, functionality, and data sources of AI-based processing transparent. The provisions on automated decision-making must be adhered to, and in the case of LLMs that communicate directly, users have a legal right to know whether they are interacting with a machine and whether the data they have entered is being processed to improve self-learning programs or used for other purposes. The use of programs that enable the falsification of faces, images, or voice messages of identifiable persons must be clearly indicated unless unlawful due to prohibitions under criminal law. A DPIA must be conducted for permitted AI-supported data processing, whereas applications that aim to undermine privacy and informational self-determination are prohibited. It is important to remember that banks offering goods or services in the European Union (EU) or European Economic Area (EEA) or monitoring the behaviour of persons in the EU/EEA also need to adhere to the General Data Protection Regulation (GDPR).

---

14–15   Article 6 FADP

16      Article 31 FADP, unless the data subject has made the personal data generally accessible without prohibiting any processing according to Article 30 para. 3 FADP

17      See 🔗 Vischer.com or 🔗 datenrecht.ch, Courageous "Hamburg Theses on Personal Reference in Large Language Models" (2024).

18      according to Article 22 FADP

19      See 🔗 Data protection impact assessment created by the Federal Data Protection and Information Commissioner. For reference, foreign data protection authorities in Germany require a DPIA for the use of artificial intelligence to process personal data to control interaction with data subjects or to evaluate personal aspects of data subjects, such as telephone call evaluation.

20      Article 21 FADP

- **Banking Secrecy:** When using GenAI in the banking sector, special consideration must be given to banking secrecy according to Article 47 of the Banking Act (BA). The disclosure of a secret entrusted to someone in his capacity in a bank is subject to criminal sanctions. Thus, Client Identifying Data (CID) should not be used by employees for prompts or training of GenAI models. When using external providers, special contractual conditions and protective measures are required.

- **Intellectual Property:** Copyright law generally grants the author the right to determine whether, when and how their work is used. Training GenAI models with copyright-protected artefacts typically requires the author's consent (unless an exception applies). Banks training their own AI models must thus consider whether their training data set contains copyright-protected material and, if so, whether the copyright holder's consent is necessary and has been obtained. The output of an AI model can also infringe intellectual property rights if copyright-protected works remain recognisable or if they contain third-party trademarks. If banks use GenAI models of third-party providers, they should carefully read the terms and conditions as they often include provisions on intellectual property and other related aspects. For example, the terms and conditions might exclude commercial use of the output from GenAI and/or reserve rights to the output for the provider.

- **Unfair Competition Act:** The Federal Act on Unfair Competition (UCA) also contains general provisions that might be relevant for GenAI applications. Generally, any conduct or business practice that is misleading or which otherwise violates the principle of good faith, such that it influences the relationship between competitors or between suppliers and customers, is unfair and unlawful.[21] This general principle is further specified in several provisions of the UCA. For example, companies (including banks) using a chatbot must avoid the misleading impression that the customer is interacting with a human if not already obvious from the context, or that an advertising picture for a product or service generated with AI is not representative.[22] In addition, companies must ensure that the statements produced by a chatbot are correct and not misleading[23]. As with other systems, they should avoid particularly aggressive sales methods using dark patterns.[24] Furthermore, companies should avoid that, for example, the output from GenAI (e.g. for marketing purposes) may lead to an unfair likelihood of confusion with competitors,[25] an unfair comparison with competitors[26] or an unfair exploitation of another person's work.[27]

- **FINMA:** Apart from the set out above FINMA Guidance 08/2024 other supervisory circulars and guidances must be taken into account when planning and implementing GenAI application. In the context of banking and GenAI especially outsourcing[28] and operational risk management[29] are expected to play an important role.

---

21–22   Article 2 UCA

23–26   Article 3 para. 1 lit. b, h, d, e UCA

27       Article 5 lit. c UCA

28       See 🔗 FINMA Circular 2018/3

29       See 🔗 FINMA Circular 2023/1

**International considerations**

- **EU AI Act:** The Artificial Intelligence Act of the European Union (EU AI Act) entered into force on August 1st 2024 and established a common regulatory and legal framework for AI. The EU AI Act classifies AI applications based on their risk of causing harm (unacceptable, high, limited and minimal) and contains specific provisions on general-purpose AI. The EU AI Act contains a broad extraterritorial scope and thus might also be applicable to Swiss banks, e.g. if they are providers or deployers of AI systems whose output is used in the EU[30]. If this is the case, the role and requirements for each application should be assessed.[31]

While no regulatory framework rivals the extraterritorial reach of the EU AI Act, Swiss banks with international activities must carefully navigate the evolving AI regulations in each country where they operate. Successfully deploying GenAI across a bank requires navigating this complex global regulatory landscape, ensuring compliance across jurisdictions while leveraging AI's potential.

## Details on risk management

Banks must establish a rigorous risk management framework to mitigate the potential risks associated with GenAI deployment. This includes but is not limited to:

- **Compliance risks:** Banks should implement a compliance monitoring system that tracks how AI systems adhere to legal and regulatory requirements. Banks should be able to demonstrate compliance in audits and regulatory reviews, especially in areas like AML, KYC (Know Your Customer), and fraud prevention.

- **Bias:** Banks should proactively address potential biases in AI systems, particularly in critical areas like credit scoring, loan approvals, and fraud detection. Bias detection and correction mechanisms should be embedded into the model development process to prevent discrimination. Additionally, algorithm transparency is essential for fostering trust and accountability. While algorithms are often proprietary and considered confidential business information, banks can demonstrate their commitment to ethical AI by engaging independent auditing companies to review their explainable AI (XAI) frameworks. These audits ensure fairness while maintaining confidentiality and signalling prudence to stakeholders, including clients, the public, and civil society organisations.

- **Customer protection:** Safeguarding customers from potential negative impacts of AI is paramount. Banks should ensure that human oversight mechanisms are in place for decisions with legal consequences or adverse effects. Importantly, they must offer clients, employees, or other stakeholders a process to appeal automated decisions they believe to be faulty. While appealing decisions such as credit denials

---

30    Article 2 para. 1 lit. c EU AI Act

31    Several law firms in Switzerland offer free assessment tools and checklists to assess potential applicability of the EU AI Act in specific cases.

may incur costs, it also provides an opportunity to improve AI systems by retraining models on identified false negatives. Offering such mechanisms not only protects stakeholders but also fosters a sense of agency, reduces reputational risks, and enhances trust in the bank's AI systems.

- **Model performance and reliability risks:** Banks should establish a model risk management framework that assesses AI models' performance, reliability, and impact. This includes regular stress testing, validation, and recalibration of models based on changing market conditions and customer behaviour.

- **Operational risks:** Banks must ensure that AI models are reliable and function as expected within operational environments. This includes testing models under various scenarios to avoid system failures that could disrupt services or lead to incorrect decisions.

- **Reputational risks:** Address reputational risks stemming from the misuse or misinterpretation of AI technologies. Transparent communication with stakeholders regarding the bank's AI practices can help build trust and credibility.

- **Cyber and information security risks:** Implement robust measures to protect GenAI systems and data from threats such as prompt injection attacks, data breaches, and unauthorized access. Securing inputs and outputs is essential to prevent manipulation or leakage of sensitive information.[32]

---

32    See for example 🔗 Vischer, Teil 6: Die andere Seite der Medaille: Wo wir KI vor Angreifern schützen müssen (2024)

## What is the view of the Swiss Financial Market Authority FINMA?

The FINMA Guidance 08/2024[33] focuses on governance and risk management for financial institutions using artificial intelligence (AI). It highlights the need for effective governance, risk classification, data quality, testing, detailed documentation, and ongoing monitoring to manage AI-related risks. The guidance emphasises the importance of explainability and independent review of AI applications to ensure transparency and accountability. FINMA also notes the challenges posed by decentralized AI development and the reliance on third-party providers.

In general, FINMA has identified four significant challenges in AI usage and expects the financial industry to address them accordingly.

- **Governance and Responsibility:** AI can autonomously make decisions or significantly influence them, complicating the control and responsibility for these actions also due to reduced transparency. This increases the risk of unnoticed errors and blurred responsibilities, especially in complex, organisation-wide processes without sufficient in-house expertise. This applies to GenAI in particular. Clear roles and responsibilities must be established, and robust risk management processes must be in place. Decision-making responsibility cannot be delegated to AI or third parties and all departments involved must have adequate AI expertise.

- **Robustness and Reliability:** AI's reliance on large datasets introduces risks from poor data quality (e.g., non-representative data). AI models self-optimize, which can lead to incorrect developments known as "drift." Increased AI usage, outsourcing, and cloud services also raise IT security risks. When developing, training, and using AI, institutions must ensure AI results are accurate, robust, and reliable by critically evaluating data, models, and outcomes.

- **Transparency and Explainability:** AI applications often involve numerous parameters and complex models, making it difficult to understand how individual parameters affect outcomes. Without a clear understanding, AI-based decisions may not be verifiable or explainable, complicating oversight by institutions, auditors, or regulators. Customers must be informed about AI usage to fully assess risks. Institutions must ensure AI results are explainable, and its use is transparent, appropriate for the recipient, relevant, and integrated into processes.

- **Non-Discrimination:** AI applications using personal data for risk assessment (e.g., setting tariffs, lending) or developing customer-specific services may yield biased or incorrect results if data on certain groups is insufficient. This can lead to unintended and unjustified discrimination, posing legal and reputational risks. Firms must ensure AI applications do not result in unjustified discrimination.

---

33     🔗 FINMA, FINMA Guidance 08/2024 (2024)

## 3.3    Technological considerations

Resource Planning | Pilot or MVP | Infrastructure and Tools | Performance Monitoring

A modern, scalable IT infrastructure is the backbone of any GenAI implementation. Banks need to ensure their technology stack can support the complexities of AI while integrating with existing systems. Key tasks include:

- **Resource Planning:** At this point, resources necessary for the GenAI implementation should be allocated (e.g., budgets, IT capacity, external consultants). Project owners should identify key roles and responsibilities, including a Center of Expertise for GenAI.

- **Launch Pilot or MVP:** Deploying a pilot or Minimum Viable Product (MVP) allows banks to validate GenAI solutions on a small scale before full implementation. This step helps assess technical feasibility, regulatory compliance, and user acceptance while identifying areas for improvement. Pilots provide essential insights into integration and performance, ensuring readiness for broader deployment with minimised risks.

- **IT Infrastructure and Tools:** When deploying GenAI solutions, selecting the proper infrastructure is critical to ensure performance, security, compliance, and scalability. The infrastructure underpins how AI models are trained, deployed, and maintained.

- **Performance Monitoring:** This task involves implementing MLOps (Machine Learning Operations) for ongoing model monitoring and updates. Success can be measured using KPIs like cost savings, customer satisfaction, and operational efficiency.

## Details on IT Infrastructure and Tools

One of the key decisions a bank must take is whether to deploy GenAI solutions on cloud infrastructure or on-premises systems. Each option has distinct advantages, challenges, and trade-offs, particularly in terms of cost, security, scalability, and compliance.[34] When planning the infrastructure for GenAI in banking, the following factors should be considered:

- **Cost:** Infrastructure costs can be significant for GenAI deployments, especially when considering hardware for training models, storage, networking, and maintenance. Banks need to balance infrastructure costs with their business objectives, ensuring they get maximum value.

---

34    See 🔗 [SBA, Cloud Guidelines (2020)](#)

- **Latency, Performance and Scalability:** For GenAI models to be usable, they require low-latency, high-performance infrastructure. This is particularly important for applications that require real-time data processing. The chosen infrastructure should be able to scale as the demands on AI models grow, whether it's expanding computing resources or handling large datasets.

- **Cybersecurity:** Banks in Switzerland are well regulated, requiring them to comply with laws and regulations like data privacy, banking secrecy and AML. Given the sensitive nature of financial data, cybersecurity should be a top priority. Banks must implement security protocols that protect AI systems from threats such as data breaches, adversarial attacks, malicious prompt injection or data poisoning. In this context, advanced encryption, sound access management, and regular security audits are essential.

- **Data Management and Integration:** GenAI systems rely on vast amounts of structured and unstructured data. Banks must ensure robust data management practices, including secure data storage, real-time data processing, and integration with legacy systems that may not easily integrate with advanced AI models. Creating a unified data platform can streamline the use of AI across departments.

- **Data Residency and Sovereignty:** Under data protection laws customer data might be required to remain within specific geographic locations. Infrastructure should support data residency requirements, ensuring data is stored and processed in compliant jurisdictions.[35]

## Details on Product and Model Considerations

The range of GenAI products and models is very broad and there are many different flavours. Whether commercial or non-commercial, open-source or proprietary, the choice of model is very much use case specific and depends on the bank's objectives, technical capacity, and legal and regulatory considerations. In general, most models can be assigned to two categories: **commercial (closed source)** or **non-commercial (open source).**

- **Commercial models** such as OpenAI's GPT4o, Anthropic's Claude 3.5, or Google's Gemini are all-rounders that can be used for various applications. Some commercial models are readily available through an API and offer plenty of features and customisation options but may lack the level of customisation banks are accustomed to. However, these models often come with their own easy-to-use interfaces, which makes them relatively straightforward to introduce to employees. Commercial models are generally made to fit the market, available as a service, and thus, reasonably easy to obtain.

- **Open-source models** such as Meta's Llama2, Mistral, or Stable Diffusion are generic general-purpose models. Because parts of the source code are open, they offer more customisation and flexibility than commercial models. Still, they are more costly to operate because of increased responsibility for security and maintenance. With open-source models, banks can tinker with the code, make the model their own, and decide whether to host it on private infrastructure or in the public cloud. It is essential to mention that most models marketed as open source are not genuinely open. They might have restrictive licensing, offer no insight into the data used to train the model, lack training code, or release only pre-trained weights without the exact compute setup.

---

35      See 🔗 SBA, Cloud Guidelines (2020)

Regardless of model or product choice, it is essential to mention that the GenAI market is very young. This space will likely experience market consolidation of products and models stemming from provider exits, acquisitions, and product discontinuations. In that sense, banks must proactively manage third-party risks by continuously assessing suppliers and establishing contingency plans. Ensuring resilience against disruptions is crucial to maintaining the quality, security, and governance of generative AI implementations.

## Deep Dive: How can the GenAI model be tailored to the needs of a bank?

A generic LLM, used in inference, has immense knowledge gathered from publicly available information, but banks typically want to go a step further and engineer a GenAI application to fit specific banking needs. Essentially, there are three ways to fit generic models to specific financial tasks, datasets, and environments. This is crucial for maximising accuracy, relevance, and compliance.

1) **Domain adaptation** is a technique employed to tailor a pre-trained model for improved accuracy in a specific domain (e.g., banking) by imparting its jargon without the need for complete model retraining. It helps the model transfer knowledge from the general domain (source domain) to the specific banking domain (target domain).

2) **Fine-tuning** involves feeding a generic, out-of-the-box LLM on proprietary information for the LLM to use this knowledge when generating outputs. For example, if an employee asked a large language model to create a customer report on portfolio performance using an out-of-the-box LLM, the AI would produce a report drawing from its generic knowledge of how customer reports on portfolio performance would look like. While the results might be feasible, they might not fulfil the bank's standards, internal processes, and corporate identity. Fine-tuning an LLM with existing customer reports would enable the LLM to mimic details, tone of voice, etc. and thus align the output better to a bank's quality standard and the habits of the customers.

3) Next to fine-tuning, there is a more straightforward approach called **Retrieval Augmented Generation (RAG).** RAG is a technique that can leverage traditional search techniques to add to a user's prompt the information most relevant to that question from a pre-defined knowledge base. The LLM can then process both the question and the additional information to generate an answer. This has proven an effective way to add company-specific knowledge to generic LLM chatbots, which is why RAG has become the most popular method to develop chatbots that are retrieval engines for a company's data. While RAG is a powerful technique, building and maintaining a robust document pipeline over time requires significant effort. However, it is essential to ensure quality answers at scale.[36]

---

36      🔗 Glean, A complete guide to retrieval augmented generation vs fine-tuning (2025)

In general, the more training a GenAI application receives, the higher the quality of responses and the less likely an LLM will be to hallucinate. RAG is one of the main ways to reduce the risk of hallucinations in a way that it can exclusively draw the LLM's attention to a pre-defined knowledge base and nothing else. At the same time, using RAG provides references that can be easily consulted by its user, if the output needs to be fact-checked. Beyond hallucinations, LLM's reasoning capabilities are also nascent. It is not unusual to find faulty logic flows in responses or have LLMs retrieve the wrong information. Various techniques can be employed to limit that risk. For instance, few-shot learning involves inputting appropriate examples of question-answer pairs directly into the prompt. Alternatively, there is the option to implement agentic-based workflows, where an LLM is explicitly asked to divide its work into specific tasks first before calling the required tools to fulfil a given task. For instance, this can be internal or external research. Once all sub-tasks in the agentic-based workflow are completed, the AI model will compile newly generated content into a final answer. This is paving the way for more complex use cases while maintaining high output quality.

# 4    What is next? – Agentic AI and its Potential in Banking

As the banking industry explores the transformative potential of GenAI, a new frontier is emerging: Agentic AI. Unlike traditional AI systems designed to assist with specific, predefined tasks, agentic AI systems can reason, plan, and act autonomously within specified parameters. These systems can integrate various tools, adapt to changing contexts, and iterate on their own outputs, making them an exciting and challenging advancement in AI development.

However, it is essential to distinguish between agents and workflows. **Workflows** are systems where large language models (LLMs) and tools are orchestrated through predefined code paths, offering predictability and consistency. **Agents,** on the other hand, dynamically direct their own processes and tool usage, maintaining control over how tasks are accomplished. This distinction helps banks decide when to implement flexible, model-driven decision-making or stick to structured, predictable solutions.

When implementing agentic AI, keeping things simple is crucial. For tasks with predictable outcomes, workflows provide an ideal solution. They ensure consistency and efficiency in well-defined processes. However, for open-ended problems that demand flexibility and adaptability, such as handling complex customer queries or creating tailored investment strategies, agentic AI systems are better suited.

Furthermore, to implement agentic AI effectively, banks can adopt practical procedures such as **prompt chaining** and **parallelisation**. Prompt chaining decomposes a task into a sequence of steps, where each LLM call builds on the output of the previous one. For instance, a bank might use prompt chaining to first generate a marketing proposal, validate it against compliance criteria, and then translate it into multiple

languages. This method increases accuracy by simplifying each step. Parallelisation, on the other hand, splits tasks into independent subtasks that are executed simultaneously. For example, analysing a client's portfolio and generating market insights before fusing both pieces of information into one single output can occur in parallel, significantly reducing processing time. Additionally, running multiple versions of the same task can offer diverse perspectives or improve confidence through consensus.

**"Unlike traditional AI systems designed to assist with specific, predefined tasks, agentic AI systems can reason, plan, and act autonomously within specified parameters."**

Transparency is paramount in these systems, particularly in a highly regulated sector like banking. Financial institutions must ensure that agents' reasoning and processes are transparent, providing stakeholders with the confidence that decisions are made responsibly and within compliance boundaries. Furthermore, extensive testing in sandbox environments is essential to mitigate risks. Controlled environments allow banks to identify potential issues before deployment. Additionally, implementing guardrails, such as limiting agent autonomy or defining clear stopping conditions, can prevent errors and thus enhance trust in these systems.

While the benefits of agentic AI are substantial, challenges such as data governance, security, and explainability must be addressed. Banks need robust governance frameworks, scalable IT infrastructure, and a skilled workforce capable of overseeing and collaborating with these systems. Pilot projects in controlled environments allow institutions to explore capabilities while mitigating risks. By iteratively refining their strategies, banks can scale these solutions effectively.

Agentic AI has the potential to redefine the banking landscape, transforming operations and elevating customer engagement to unprecedented levels. Starting small, prioritising simplicity, and maintaining transparency will enable financial institutions to harness this transformative technology responsibly. While the field is still evolving, agentic AI's trajectory signals a future where intelligent, autonomous systems become integral to financial success.

## Practical examples: what could be the impact of agentic AI on banking?

A practical example of agentic AI in banking is customer support. Consider an online banking chatbot enhanced with external tools. Such a system could retrieve a client's payment data, analyse spending patterns, and generate personalised savings plans. By integrating these capabilities, the chatbot could not only answer queries but also perform programmatic actions, such as issuing refunds or updating account details. These dynamic interactions elevate customer engagement and demonstrate the potential of agentic AI in delivering tailored financial solutions.

But: not only banks make the most of powerful AI applications. Most likely, customers will soon have access to AI as powerful as that of banks, allowing them to autonomously evaluate financial offers, manage transactions, and negotiate financial terms. The consequence of this shift might be that customers would no longer interact with banks directly. Instead, their AI bots would negotiate financial services on their behalf. Having permissioned access to bank accounts via open banking frameworks, these bots would be capable of optimising financial decisions in real-time, selecting the best products based on price and terms rather than brand loyalty or user experience. In that scenario, competition would shift purely to pricing and service efficiency, putting trust at the core of the interactions – banks' most considerable advantage and something that fully AI-based services will struggle to establish.[37]

---

[37]     Birch, D., Rutter, K. Where are the customers Where are the customers' bots? The AI paradigm shift in retail banking. (2023)

# 5    Conclusion

**GenAI is a transformative technology with the potential to significantly enhance banking operations.** However, deploying this technology in a highly regulated industry demands careful consideration and balance between innovation, compliance, governance, and risk management.

**Key to this balance is understanding both the capabilities and limitations of GenAI.** Techniques like Retrieval Augmented Generation (RAG) and advanced prompting strategies, such as few-shot learning and agentic-based workflows, play a crucial role in improving the accuracy and quality of outputs while minimizing risks like hallucinations and faulty logic. Despite these advancements, deploying and maintaining such systems requires robust IT infrastructure, clear governance, and comprehensive training for users to extract the most value from the technology while addressing adoption challenges.

**A tailored implementation strategy is essential, as there is no one-size-fits-all approach.** Banks must align GenAI adoption with their specific business models, regulatory requirements, and technological maturity. Successful integration should prioritize enhancing employee productivity and human expertise. A structured framework, differentiating tasks across various phases, provides valuable guidance, but banks must adapt it to their unique circumstances.

**Grassroots adoption by employees is a key driver, but management should provide a structured oversight.** Many employees are already experimenting with GenAI tools like ChatGPT in their daily workflows, highlighting the technology's potential. While this organic adoption fosters innovation, it also underscores the need for a formalized strategy. Banks should embrace an iterative approach where employee experiences inform leadership decisions, ensuring GenAI is implemented securely and in compliance with industry requirements and standards.

**Emerging trends, such as agentic AI, require further exploration and monitoring.** While agentic AI introduces more autonomous and adaptive workflows, its full implications – particularly regarding transparency, accountability, and regulatory compliance – are yet to be further analysed and fully understood. Banks must approach these advancements with the usual deliberateness, ensuring responsible implementation that balances innovation with oversight.

**This paper serves as an initial outline and point of reference rather than a definitive roadmap.** While it highlights key considerations for GenAI adoption in banking, most vital aspect of AI and GenAI use, namely assessing the risks and defining appropriate technical and organisational measures to deal with them, remains the responsibility of each individual institution.

**Ultimately, the successful integration of GenAI requires acting across multiple dimensions.** Aligning AI initiatives with business goals, fostering a culture of innovation, ensuring regulatory compliance, and implementing rigorous risk controls. By addressing these dimensions comprehensively, banks can unlock the transformative power of GenAI, enabling them to remain competitive in an ever-evolving financial landscape.

# Glossary

**Automated Decision-Making:** AI-driven decisions made without direct human intervention, subject to legal and compliance constraints.

**Data Governance:** The management of data quality, security, and compliance to ensure AI systems function ethically and legally.

**Data Residency:** Regulations requiring customer data to be stored and processed within specific jurisdictions to comply with data protection laws.

**Explainability:** The ability to understand and interpret AI decision-making processes, critical for regulatory compliance and customer trust.

**Fine-Tuning:** The process of adapting a pre-trained AI model to a specific banking use case by training it on proprietary data to improve accuracy and relevance.

**Latency:** The time delay in AI processing, which is critical in real-time banking applications like fraud detection.

**Model Drift:** The degradation of AI model accuracy over time due to changes in data patterns, requiring continuous monitoring and retraining.

**MLOps (Machine Learning Operations):** The practice of deploying, monitoring, and managing AI models in production environments to ensure performance and compliance.

**Parallelization:** An AI technique that splits tasks into independent components for simultaneous execution, improving efficiency.

**Prompt Chaining:** A method where AI responses are generated step-by-step, ensuring accuracy and logical consistency.

**Retrieval-Augmented Generation (RAG):** An AI technique that combines retrieval-based search with generative models to improve accuracy and reduce hallucinations.

**Tokenisation (AI Context):** The process of breaking down text into smaller units (tokens) for efficient AI processing.

**Workflows (AI Workflows):** Predefined sequences of AI-driven tasks that ensure consistency in banking operations, distinct from fully autonomous agentic AI.

# Edited by

**Andrea Luca Aerni,** Policy Advisor Digital Finance, SBA
**Richard Hess,** Head of Digital Finance, SBA
**Tomas Marik,** Intern Digital Finance, SBA
**Panagiotis Psomas,** Intern Digital Finance, SBA

# Experts

**David Adametz,** Senior Data Scientist, Center of Excellence for AI, UBS Switzerland AG
**Steve Blanchet,** Head of Group Technology Strategy and Innovation, Banque Pictet & Cie SA
**Frédéric Dommart,** Chief Data Officer, Bank Lombard Odier & Co Ltd
**Julinda Gllavata,** former Head Data Analytics & AI, SIX Group Ltd
**Simon Gomez,** Head Data and Innovation, LGT Private Banking
**Stefan Jeker,** Head Innovation Management, Raiffeisen Switzerland Cooperative
**Jürgen Petry,** Open Finance Lead, Raiffeisen Switzerland Cooperative
**Matthias Plattner,** Head Channels & Digital Services, Bank Julius Baer & Co. Ltd.
**Prof. Dr. iur. Cornelia Stengel,** Co-Director, Swiss FinTech Innovations (SFTI)

## About the Swiss Bankers Association (SBA)

The Swiss Bankers Association (SBA) is the Swiss financial sector's leading industry organisation and represents the interests of some 270 member institutions. Founded in 1912, it ensures optimal framework conditions for a competitive and innovative Swiss banking industry. It promotes dialogue with politicians and authorities, drives vital topics such as sustainable finance and digital currencies, and supports education and professional development in the industry. As a knowledge center, it is dedicated to the sustainable development of the banking industry.

swissbanking.ch