



Financial Crime Fighting and Risk Control in the Age of Generative AI

Transformational Improvements of Scale,
Effectiveness, and Efficiency

Peter Cousins, Chief Technology Officer, WorkFusion

Reducing Fraud, Risk, and Illegal Money Flows – A Worthy Challenge

Reducing crime and risk through financial system monitoring and controls has been a mixed success. Despite the large investments made by the public and private sectors, it is never enough to keep pace with the increasing volume and sophistication of bad actors. LexisNexis estimated the global cost of compliance was nearly \$275B in 2022, with 60% of that tied to labor (direct and outsourced).¹ The labor-intensive processes of today restrict the scope of these programs to the point that they are mostly focused on minimum viable regulatory compliance. We must move past our current practices through recent technological advancements to fundamentally move the needle on risk and the greater mission of social responsibility.

The fundamental blueprint for this new era is the same as it ever was:

1. **Monitoring** a sprawling landscape of data for events that require attention.
2. **Reconciling** these events with entities of interest.
3. **Enriching** the context using data providers, web data, and enterprise data.
4. **Requesting** required information from customers for validation or explanation.
5. **Reviewing** unstructured source material to extract salient data points and meaning.
6. **Curating** and computing metrics on the broader data set based upon domain expertise.
7. **Eliminating** events that represent immaterial changes from previous investigations.

8. **Analyzing** events by considering previous behavior, expectations, and contextual data.
9. **Finding** anomalies to prevent new threats from going undetected.
10. **Creating** a dossier from all this information, making decisions, and explaining reasoning.
11. **Acting** quickly and directly when problems are discovered.
12. **Ensuring** correct results through quality control and continuous improvement.

These things are each a substantial amount of labor to perform in the traditional way. The cognitive load, sheer amount of data, and the frequency of cases closed without incident creates a constant risk of inconsistent execution and analysis. It is an incongruous mixture of mind-numbing routine review and the fear of the consequences of missing critical incidents.

Analysts work under time pressure, partially because of the labor expense driving pressure to reduce unit costs, but also because of revenue impact. Backlogs create customer satisfaction problems, risk customer attrition, and delay revenue producing activities.

The work is detailed and requires significant training to perform correctly. The difficulty in retaining employees makes it very difficult to make meaningful gains in the overall strength of the team. The workload can also be unpredictable, making it difficult to predict the team size required to handle peak workloads without unnecessary overhead during non-peak times.



This generalized blueprint and problem statements can be applied to virtually all aspects of monitoring and investigations. However, Know Your Customer (KYC) processes have been some of the most expensive and challenging to improve.

It has been stubbornly difficult to automate these tasks due to the high variability, unstructured information, research, outreach, and judgment – until now.

pKYC: Myths and Legend

Perpetual KYC (or pKYC) has been discussed for years. If done right, it has the potential to greatly improve on the status quo. PwC estimates upwards of 60-80% cost savings for effective pKYC implementations.²

In traditional KYC, a periodic refresh of the KYC information and process is performed on a schedule according to the risk categorization of a customer. A low-risk customer might be refreshed every three years, whereas a high-risk customer would be refreshed annually. This is both too much and not enough: often the refresh is simply reiterating the information already known, and problems can often occur in between the intervals.

In pKYC, the idea is to monitor and react to events whenever they occur. In this way important changes are investigated promptly rather than waiting for a fixed refresh date, and clients with no changes are not unnecessarily refreshed simply because time has passed.

The problem with pKYC implementations so far has been that they have focused on acquiring the events without automating the actual handling of these events. This creates an order of magnitude increase in the workload, so it becomes impractical to keep up and completely undermines the potential benefits. A partial solution is thus worse than useless, and these projects cannot really be considered pKYC.

For pKYC to become successful and mainstream, we must automate every aspect of the process, and only involve people when it is necessary to seek secondary opinions, escalate, or handle unknown situations. Unknown situations can include difficulty obtaining information through automated means, insufficient similar examples for training, or a risk identified that warrants human involvement.



Automation needs to focus on:

- Acquisition of events that drive pKYC, starting with data providers but moving beyond.
- Ingesting the process flows that comprise the standard operating procedure.
- Executing the process flows with a mix of automated steps and human involvement.
- Continuously converting steps that require human involvement to automated versions. Steps that are run a larger number of times drive an empirical heatmap and provide the supervised examples needed to train new models.
- For low-risk populations, automating the process completely except for escalations.
- Automating assessment/disposition of immaterial events, no matter the client risk level.
- For higher-risk scenarios, completing the dossier, analysis, narrative, and proposed disposition – subject to accelerated human review and approval.

¹ Lexis Nexis Risk Solutions: [2022 True Cost of Financial Crime Compliance Study](#)

² PwC: Perpetual KYC: [A new approach to periodic reviews](#)

How Generative AI and Large Language Models Can Help

Many of the challenges automating pKYC have their roots in the complexity of the natural language processing, structuring of the information, and working collaboratively with a person without losing the benefit of partial automation. Generative AI and Large Language models (LLMs), such as Google Bard/PaLM 2, can drive significantly higher performance in many areas.

Synopsis Generation

Summary of Document, Email, Case, Dataset, Decision

Throughout a pKYC process, synopsis generation can help drive collaboration between the automation and the people involved.

When a person intervenes in an investigation, there may be a detailed set of facts that have been gathered in a dossier, which can be summarized to quickly orient the analyst on the progression of the case.

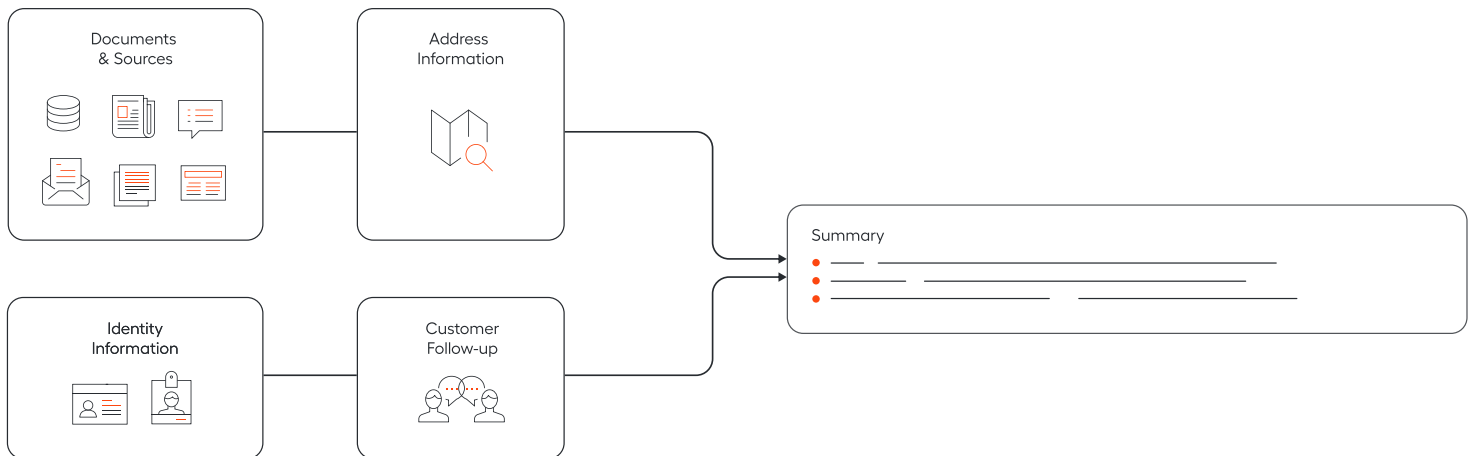
There may be many individual documents that have been gathered, and these documents can be lengthy. A synopsis of each document allows the analyst to more easily identify the ones that warrant closer inspection.

There may be customer outreach that has not been able to acquire the necessary information. The email thread with the customer can be summarized to explain what was requested, what was correctly supplied and what issues are still outstanding. It can include an explanation of why outreach has been escalated, such as non-responsiveness, customer apparent confusion, refusal to supply information, multiple incomplete or erroneous submissions, apparent frustration, or customer requested escalation.

A synopsis can also be used to elaborate on the structured information to highlight unusual factors or comparisons with historical information about the customer or their peer group. A complex set of chained ultimate beneficial ownership statements can be summarized to collapse unnecessary detail. A summary of the individual identity documents including generalizations about individuals, addresses, and roles can make digesting this information easier.

Decisions can also be complex, layered and multi-factor, so having a synopsis generated from the complete reasoning can accelerate understanding of the most important issues. This can also point to other information gathered in the dossier and include synopsis of this information in the citation explaining the reasoning.

All the necessary detail is compiled in the dossier, which can be consulted to clarify and validate the synopsis if it seems ambiguous, confusing, or implausible. This use of the technology can be a significant efficiency gain for the human collaboration, and the risk is mitigated by the human supervision in consuming the synopsis and the supporting detail.



Answering Questions – Key terms, entities, highlights, unusual elements

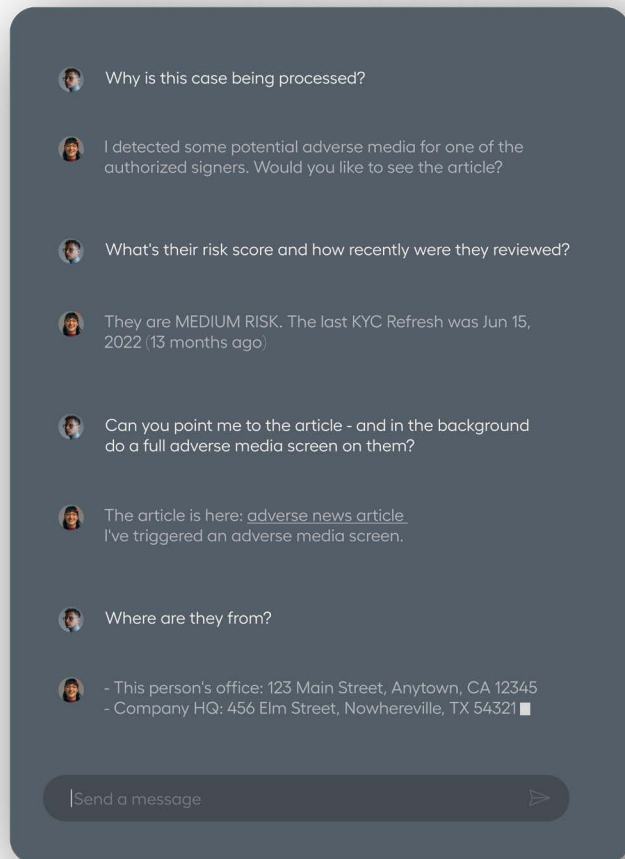
Throughout a pKYC process, using generative AI to answer questions can help drive collaboration between the automation and the people involved.

Given the large amount of information in the dossier, an efficient way to navigate the information is to ask questions which can be extracted from the complete information set, including attached documents.

Questions can include:

- Information about associated entities and their geographical information.
- Information about what parts of the process have been completed.
- Clarifications about the original event signal that triggered the process.
- Details about change detection and comparisons of the old and new information.
- Questions about availability or the nature of the additional data gathered.
- Information about customer outreach status and who from the customer or the bank has been involved thus far in the information outreach conversation.
- Information about key terms in any document, event signal, or web data gathered.
- Questions requesting characterization of unusual elements.
- Questions about sentiment of related news articles, reports, or filings.

Question answering is similar to synopsis generation, in that it is effectively a supervised activity. It is a productivity improvement, but the analyst can verify or clarify any answers that are confusing, ambiguous, or suspect. This leverages the technology while managing risk effectively.



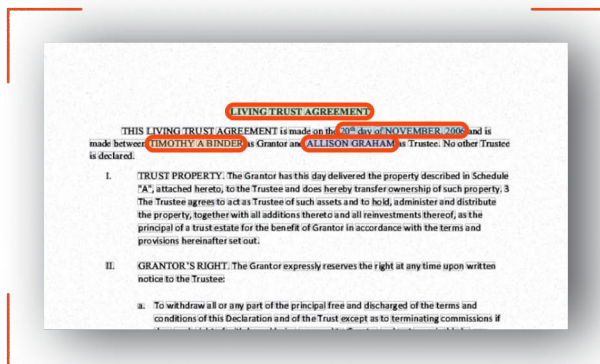
OCR and Input Corrections – Analysis of improbable text and automatic correction

Many of the data elements required during the execution of the pKYC process can have subtle errors that require correction for successful automation. These errors, if left unresolved could fail to satisfy the processing criteria leading to unnecessary exception processing.

The problems can come from manual input typographical errors, scanning a document using OCR that contains visual noise preventing perfect recognition, or handwriting that is intrinsically less reliably recognized. Text can even come from transcriptions of voice input, which is notoriously riddled with word level recognition errors.

Traditionally statistical approaches are used to predict errors or select most likely best interpretation. For input text, corrections might compare Damerau-Levenshtein distance in determining likely typographical errors, or abbreviation alias maps to correct abbreviations to formal text. For OCR text, most likely text is selected based upon recognition confidence, but the highest probability selection is followed by less confident alternatives that may be the actual correct value.

With new language models, semantic based predictions on likely text can be compared with the traditional techniques to determine better automatic corrections. While an OCR engine might see a name and interpret the most likely text as 1imothy, contextual analysis would recognize it an improbable name and provide Timothy as the top correction. In such cases, a human analyst may be requested to validate the correction, but the time to resolve is greatly reduced by increasing the likelihood of correct values being selected.



Document Extraction – For automated extraction QC, or HITL acceleration

Intelligent Document Processing is a key capability required during pKYC. Even given high fidelity digital ingestion of the document text, there is substantial processing required to extract structured data from the documents necessary to complete the pKYC process. These elements are individual fields that will be further processed and compared against systems of record or to drive rules and conditions for next steps.

Traditional methods yield good results, but when highly unstructured data is involved, there is always a risk of information extraction failure – either failure to find a data element for extraction, extraction of a subset of the correct information for the element, selection of the incorrect data element, or poor coverage of a completely uncovered example due to a long tail of document variations.

The information extraction capabilities of large language models are complementary to traditional extraction techniques, because they can be used in an ensemble style processing to compare the results from both traditional and LLM approaches. This is especially helpful because if the traditional IDP and new LLM models agree, there is an extremely high probability that the extraction is correct. If the answers diverge, reconciliation algorithms can be used as a tie breaker in many cases – for example if there were adjacent text that was missed and the inclusion is credible, or if there were multiple candidates identified a common candidate can be selected automatically.

If there are extraction failures due to apparently missing data or a long tail document example, the LLM can be used to drive a human review and correction in a greatly accelerated way, and such corrections can be used in subsequent iterations of the model to improve extraction performance. This combination of ensemble model processing with human review can drive automation rates higher than ever before while keeping risk low.

Entity Substantiation – Comparison of profile with document for Identity substantiation

Some signal events come from data sources that are not definitively matched to the entities in question. In order to automate the processing of these events, subtle cues can be used to substantiate or repudiate the match.

Traditional models are used to interpret clear statements about identity that can be compared to the profile of record to see if the entity matches, such as information about an entity, related entities (such as employer for an individual, or the CEO name related to a company), geography, age, or other identifying information. This is a key element of adverse media monitoring.

New LLM techniques can be used to perform analysis of more subtle factors, such as an ambiguous reference (for example to a location as a region instead of a city), or data elements that are related but through indirect references in the text. This can supplement traditional methods of matching and can increase automation rate, or better highlight significant passages for human review and decisions.

Document Search – Finding documents for research / analytics beyond text search

For some pKYC cases, a document repository may include a significant number of documents that are either linked explicitly to the investigation or have been linked implicitly through link analysis. Such documents can be difficult to locate or identify as the appropriate candidate document.

Traditional full text search approaches wind up being a tedious and ineffective method for finding and correlating such documents due to the high incidence of false positive or noisy hits since so many documents contain some of the information in question. It is also difficult to determine missing documents.

Using LLM in search and correlation of such documents can evaluate conditions across them. For example, finding the identity documents for officers that will expire this year, or listing beneficial owners from or near a geography of concern. LLM search is more tolerant of finding data with differing keywords but the same meaning. Document search can be combined with traditional techniques to perform high fidelity data extraction on the resulting documents, and if the

structured data comparison matches, the result can be automated. If not, it can be a tool for helping an analyst more quickly find the information than with plain text search techniques.

Document Classification – Document type analysis (supervised)

Document type classification is a key part of evaluating documents that have been supplied by the customer or relationship manager to make sure it was the type of information requested. This helps make sure the process is ready to run and no further outreach is required.

Traditional approaches for defining document classification models are efficient, as is the evaluation of the classification at run time. However, there may be instances where a document cannot be classified due to long tail document types or a yet to be seen new variation of a previously seen document type.

In cases where an unknown document type is encountered, a fallback to LLM techniques for attempting a similar classification may succeed due to deeper analysis of the language. For example, a classifier may not recognize a new document as a kind of trust agreement, because it lacks key terms always present in previous trained examples – but the LLM might recognize it anyway due to more subtle correlations in the language between this and previous trust agreements. In such a case, the trust agreement extraction model could be run and produce all required data once it is identified as a trust agreement, and this would validate the accuracy of the classification. Whether full automation or additional analyst review, the net automation rate gain can be significant.

Document Clustering – Analysis of types of documents for process model coverage analysis

At various stages in the pKYC process, there may be variable documents that will be supplied. Each type of document would need to have examples trained for classification, and an extraction model created, so that appropriate decisions or actions can follow.



One challenge in building out the processing model is ensuring complete coverage of conditions. While non-covered cases can be set aside and trigger manual intervention by an analyst, it is desirable to know from historical examples that an uncovered cases exists. This then allows decision making based upon data volumes to proactively design the process to handle these cases.

LLM document clustering can be used to create a map of the likely document types and compare to the process model to determine which clusters have coverage in classification models, extraction models, and next steps. This is part of the design activity and is executed under designer supervision.

Document Correlation – Association of documents with outcomes / making predictions

During the pKYC process, a number of documents are considered. A document might be a possible adverse media article and it will be analyzed to determine identity, focality, and seriousness of an issue. Another document might be an SEC filing where information is extracted. These techniques are used to automatically process such documents as much as possible.

However, analysts also might flag an issue that is concerning during a manual review. For example, even though the entity in question is the victim in a court case, which is typically clears the adverse media article, the fact they are involved in the case at all might lead to further questions. This can be handled by relaxing the focality constraint on the media analysis.

Another way to influence the disposition of this and similar potential issues is to label documents from historical and ongoing analysis. Associating the labels of false, true, and suspect to these documents can be used to make similar predictions about future documents. This can capture emerging or unexpected subtleties and thus make the review more like an expert analyst review. It can be a form of layered protection to be used in conjunction with the traditional approaches to document evaluation.

In addition, all documents can be labeled based upon historical outcomes to indicate whether a document was associated with an issue of concern or not. This can provide a similar data point on even routine documents to identify subtle points of interest. In such cases it can trigger analyst manual review with a citation of the document correlated with negative historical outcomes, but the analyst would make the final decision.

Conversational Interactions – Training / Process Definition / Feedback

As most of the previous examples have illustrated, the pKYC process is dynamic and needs to adapt and learn in cooperation with analysts. Making the evolution a natural collaborative process is essential or it becomes an IT process instead of a subject matter expert driven one. Conversational interactions driven by LLMs are a key enabling strategy.

Adding additional labeling information can be done conversationally, which can be used to train various models that require the labels for classification or correlation.

Processes will start small, so that they are simple, and the first step is ingesting standard and procedures to create the rough process model. Corrections can be given that will tune the process when the ingestion is incomplete. Specifying such corrections using natural language is the best way to collect feedback at runtime from business users.

Collecting feedback is also essential to specify reasons for clearing a false positive or correcting the scope of an extraction or classification. Even the best rules builder is better if the rule can be specified or corrected with natural language.

These conversational interactions are inherently supervised and use generative capabilities to produce what might be otherwise an uncomfortable formal specification for an analyst.



Summary of How LLMs impact pKYC

LLM Area	pKYC Use	pKYC Impact
Synopsis Generation	A summarized dossier of docs, emails, case details, datasets, and decisions can help drive collaboration between the automation and the people involved.	All necessary case detail compiled in a single source to consult in investigations, an efficiency gain for the human collaboration, and the risk is mitigated by the human supervision in consuming the synopsis and the supporting detail.
Answering Questions	Given the large amount of information in a case, an efficient way to navigate the information is to ask questions which can be derived from the complete information set, including attached documents.	A productivity improvement allowing the analyst to easily verify or clarify any answers that are confusing, ambiguous, or suspect. This leverages the technology while managing risk effectively.
OCR and Input Corrections	Many of the data elements required can have subtle errors that require correction, which can come from manual input error, visual noise preventing perfect OCR recognition, difficult handwriting to recognize, or transcriptions of voice often containing recognition errors.	Semantic based predictions on likely text can be compared with the traditional techniques to determine better automatic corrections, enabling a human analyst to easily validate any corrections with minimal time-to-resolution by increasing the likelihood of correct values being selected.
Document Extraction	Key to the process is collecting specific data points from the available documentation, which can then be further processed and compared against systems of record or to drive rules and conditions for next steps.	Extraction failures due to apparently missing data or a long tail document types can be minimized by flagging a need for human review and correction in an accelerated way, with corrections leveraged in subsequent iterations models with improved extraction performance, thus driving higher automation rates while keeping risk low.
Entity Substantiation	Customer profiles can be compared with provided documentation for validating or invalidating identity, especially subtle factors such as comparing regions rather than specific cities	Being able to compare items in ways that are semantically based or by following indirect references enables higher automation rates or highlighting specific needs for human intervention.
Document Search	Documents for research and analysis can be located via means beyond basic text search, such as finding identity documentation about to expire or listing beneficial owners from or near geographies of concern.	By combining with other automation techniques such as data extraction, this could increase automation of some cases or help analysts more quickly find the information they need.
Document Classification	The type of document or data under review can be categorized, often for the purposes of further automation. For example, by being able to better identify which documents are Trust Agreements, more of the data associated with Trust Agreements can be captured via automation.	Net automation rates will increase, not only for the classification of documents, emails, and other data, but also in the downstream activities that benefit from more accurate and complete classification upstream.
Document Clustering	Analyzing the types of documents involved in historic processing helps determine the full scope of classification, extraction, and next steps required.	More comprehensive automation can be defined and setup, increasing the scope of what is subject to automation and what will remain manual in the near-term and future.
Document Correlation	As a number of documents are considered as part of the process, it is useful to be able to further characterize source data involved in investigations.	By better labeling data for additional content and context clues, predictions for things such as what is concerning or not concerning help analyst make stronger final decisions.
Conversational Interactions	Providing additional training, process definition, and feedback correlating to relevant natural language.	An ongoing feedback loop, this avoids the need for analysts to provide an upfront, formal specification and instead offer a more comfortable method for improving the performance of automation.



The Complete Picture

The elements of a successful implementation require more than just LLMs. A holistic approach to automation and AI is required. Our pKYC approach uses the following no-code elements:

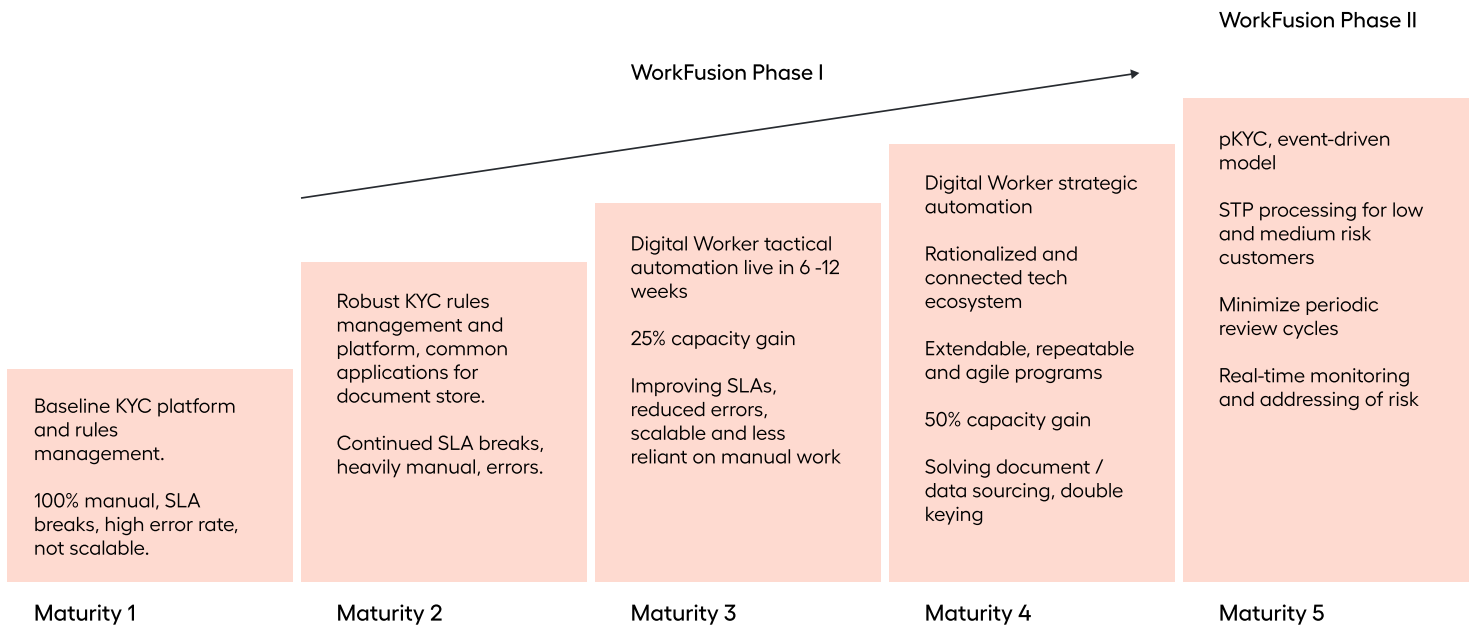
- Digital Workers who collaborate with their human colleagues, working assigned cases, and writing notes in the case management system. When cases are completed, human colleagues may perform completion or review as uncertainty or risk warrants.
- Human in the Loop tasks allow for people to be involved either during the process or as a final review. The interface is designed and implemented without UI developers.
- Working the cases means ingesting and following the standard operating model to run a pKYC process. Initially, the workflows may be more manually intensive, but ultimately all the common decisions and actions provide training examples for AI models.
- AI Models can be produced without coding but also without data science skills. AI models extract data from documents, classify data or documents, and make decisions.
- Automation covers collecting and maintaining data from the existing systems of record, complex documents, and supplemental data sources. Final actions include writing justifications for these decisions, understanding when escalation is required, or updating systems as necessary.
- Connectors for sourcing input events from any possible source without writing code. These connectors can use data providers/aggregators, web search, packaged financial services applications, cloud services, proprietary enterprise systems and middleware, documents, UI/RPA, or databases/BigQuery.
- Process modeling tools and runtime to define the steps necessary for the desired outcome. Each step can process data and make decisions. Data processing can include data enrichment, transformation, correlation, curation, correction, and change detection. In the end, the connectors are used to syndicate data and take actions.
- Quality checks and Model Risk Management build confidence on correctness and fairness of the models, where the business stays in control of the level of automation.
- Provisioning automatically in Bank VPCs on Google Cloud, and elastically scaling environments automatically.

The key to success on these projects is to break down into phases that can help drive benefits continuously. We look at this as a journey against the “KYC Maturity Model” (see diagram on the next page). Even before full pKYC is tackled as a program, significant gains can come from the data and document automation to help save time for KYC teams. Next, capacity gains can come from KYC assist through customer outreach automation, identity document ingestion and customer due diligence document ingestion. pKYC then follows, starting with full automation of low-risk populations as the first wave. In subsequent waves, automation of any low-risk signal that can be eliminated, followed by running full pKYC with human oversight on higher risk populations.

Generative AI and Large Language Models represent a key enabling technology to bring a step change to pKYC and other efforts to reduce fraud, crime, and risk. The stacking benefit of traditional techniques with new emerging technologies make these efforts easier and scalable.



KYC Maturity Model



The benefits of doing this are not only saving time and money but reducing human error and catching more problems than ever before. McKinsey estimates that enhanced KYC delivers 30-50% faster onboarding, driven by 20-40% fewer touchpoints and 20-30% faster customer RFI responses, which results in a 10-30% increase in customer satisfaction (CSAT) and 10-15% reduction in employee attrition. Large banks can save thousands of FTEs and improve customer satisfaction, time to revenue, employee attrition, and regulatory outcomes.

With fundamental improvements in automation and efficiency, we can now screen for non-regulatory issues that improve risk and avoid associations with unsavory parties. We can avoid doing business with the entities that wage wars of aggression, commit crimes, oppress people, exploit children, or destroy the environment.

If we don't do these things, it will not only be a lost opportunity – as sophisticated bad actors and states embrace the technology they will exploit our complacency and evade the values and regulations that we strive to protect.

Special thanks to James Rochford, Global Head of KYC Services, of Deutsche Bank, for his contribution to this paper.

³ McKinsey & Company: [Five actions to build next-generation know-your-customer capabilities](#)