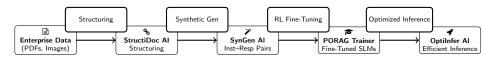
### **AI-Powered Enterprise Innovations**

S. S. Srinivas, Shivam Gupta, Akash Das, Venkataramana Runkana

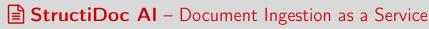
- Foundational Capability: End-to-End AI Pipeline for Document Intelligence and Model Optimization
  - StructiDoc + SynGen + PORAG + OptiInfer: A modular Al workflow that transforms unstructured enterprise documents into structured knowledge, generates synthetic instruction-response datasets, fine-tunes small language models via policy-optimized retrieval-augmented training, and deploys them with high-speed inference-time optimization.
- Domain-Specific Application I: Automation in Chemical Engineering
  - AutoChemSchematic AI: A closed-loop, physics-aware agentic framework for the automated generation and validation of chemical Process Flow Diagrams (PFDs) and Piping & Instrumentation Diagrams (PIDs) for novel chemical industrial production process.
- Domain-Specific Application II: Advanced Solutions for Competitive Advertising
  - Agentic Multimodal AI for Advertising: A framework for hyper-personalized B2B/B2C competitive advertising, leveraging multimodal market intelligence, persona simulation, and adaptive ad generation.
- Domain-Specific Application III: Transforming B2B Chemical Commerce
  - **Q** ChemConnect AI: Digitizing chemical commerce via an agentic multimodal B2B marketplace, featuring Al-driven product data structuring and actionable competitive insights.

#### % Unified AI Workflow: From Raw Enterprise Data to Optimized AI Outputs



#### Stages:

- E StructiDoc AI: Parses unstructured content into structured formats
- SynGen AI: Generates synthetic instruction—response pairs
- PORAG Trainer: Fine-tunes small language models for RAG
- **4** Optilnfer AI: Enables fast, scalable, inference-time optimization



■ Unlocking the Power of Unstructured Enterprise Data into Machine-Interpretable Knowledge

S. S. Srinivas, Shivam Gupta, Akash Das, Venkataramana Runkana

#### ■ A The Challenge: Proliferation of Unstructured Data

- Most enterprise data (e.g., scanned PDFs, contracts, invoices, spreadsheets, handwritten notes) is unstructured or semi-structured.
- Extracting structured, machine-readable data from unstructured or semi-structured data sources remains a major bottleneck for Al-powered automation, analytics, and decision-making.

#### ■ The Role of Structured Data in Enhancing LLM Performance:

- LLMs achieve factual consistency, higher accuracy, and reliability when grounded in structured, high-quality data.
- # Production Al systems demand structured inputs—yet most enterprise data is unstructured. Bridging this gap requires automated pipelines to parse, validate, and transform raw enterprise documents at scale.

#### ■ Growing Demand for Al-Ready Data

- **A**I-native enterprises increasingly require clean, structured, and explainable data:
  - Fine-tune small-scale (LMs) on domain-specific corpora for task-specific customization
  - Q Build Retrieval-Augmented Generation (RAG) systems
  - Q Power document search, summarization, and reasoning agents
- Var Proposed Solution: A document ingestion platform designed to convert complex, unstructured documents into structured, machine-readable data optimized for LLMs workflows

#### 

#### ■ A Challenges with Complex Layouts in Unstructured Documents:

 Multi-column formats, nested tables, and embedded visuals(Figures, Images, etc) often lead to misinterpretation or data loss.

#### Example:

- A research paper with 2 columns is extracted as a jumbled text stream, mixing left and right column content.
- **III** A **nested table** in an invoice is read row-by-row, losing column associations (e.g., "Quantity" vs. "Unit Price").
- Al Solution: Layout-aware multimodal models preserve structure in multi-column text, tables, and visuals.

#### 

 Traditional OCR (e.g., Tesseract OCR, PaddleOCR) extracts text but cannot infer meaning or relationships.
 Example:

#### =xampic.

- **É**\* Extracts "Apple" but cannot tell if it's **fruit** or **brand**.
- 🖹 Reads "Total: \$100" and "Due: 30/05/2024" but fails to link them as part of the same payment information.
- ◆ AI Solution: LLMs classify "Apple" by context and group invoice fields logically.

#### • A Restricted Language and Font Support:

- Struggles with non-Latin scripts, handwritten text, or stylized fonts.
   Example:
  - 🎉 A Japanese Kanji receipt is misread as random symbols.
  - Doctor's handwritten prescription is rendered as gibberish.
  - ◆ Al Solution: Multilingual transformers (e.g., VLMs such as GPT-4o) improve accuracy across scripts.

#### Dependence on Image Quality:

- Traditional OCR fails on poor scans or noisy images; requires ideal input.
   Example:
  - Blurry ID "DL 8HX" misread as "DL 8KX" by Tesseract.

#### Al Solution:

 VLMs infer noisy text(blurred/unclear text) correctly using contextual understanding.

#### ■ Security and Compliance Risks:

- Traditional cloud OCR exposes sensitive data to third-party services.
   Example:
  - Hospital records processed on external servers violate HIPAA(Health Insurance Portability and Accountability Act).

#### Al Solution:

 ■ On-device OCR systems and language-only AI models enable secure document understanding tasks—such as summarization, Q&A, etc—without relying on external servers.

#### 

#### Inflexibility:

 Robotic Process Automation(RPA) tools (e.g., UiPath, Blue Prism) are software robots (or "bots") to automate repetitive, rule-based tasks that are typically performed by humans in digital systems.

#### What RPA Tools Can Do with Documents:

- $\bullet \hspace{0.5cm} \textcircled{D} \hspace{0.5cm} \text{Open a scanned PDF} \to \text{run OCR} \to \text{extract key fields} \to \text{enter in form}$
- ullet Parse structured forms o validate values o trigger follow-up workflows
- RPA tools are effective only when document formats are rigid and predictable.
- But RPA tools breaks on slight variations in structure or layout.

#### Example:

- $\bullet \;\; \stackrel{\blacksquare}{\Longrightarrow} \; A$  new invoice template requires re-training the entire RPA pipeline.
- Al Solutions:
  - Claude 3 handles 100+ invoice formats out-of-the-box with layout-agnostic reasoning to generalize across structural variations.

#### O High Costs:

- Traditional systems require constant maintenance Example:
  - 📽 Full-Time Equivalents(FTEs) needed to correct insurance claim OCR errors
  - Al Solutions:
    - Fine-tuned SLMs automate extraction, reducing the need for manual FTE intervention.

#### Vision-Based Document Ingestion for LLM Pipelines

#### Vision-Centric Document Processing Engine:

- Uses layout-aware models to segment and classify document components—such as text blocks, tables, images, and figures—across both standard and non-standard layouts, including multi-column and nested structures
- Applies specialized extraction pipelines for each content type (e.g., figure-caption linking), preserving semantic relationships and visual hierarchy
- Reconstructs the document into a structured, LLM-ready format that retains its original meaning and context, enabling accurate downstream applications like RAG.

#### ■ ♣ Advanced Document Parsing:

- Uses multi-pass Agentic OCR with VLMs to accurately extract structured data from complex documents. Generates machine-readable formats (JSON, XML, HTML) optimized for LLM pipelines and retrieval systems.
- Supports custom schema definitions to fit domain-specific data extraction needs.

#### Deployment Flexibility:

- Supports both cloud-hosted SaaS and secure on-premises installations, ideal for regulated industries handling sensitive documents.
- Provides REST APIs and Python SDKs for both synchronous and asynchronous ingestion workflows.

#### ■ Security and Compliance:

- Enforces zero data retention—no documents are stored post-processing.
- Offers air-gapped and on-premises deployment for maximum data privacy.
- Compliant with HIPAA and SOC 2 Type 2, ensuring security and privacy controls.

#### ☑ 1. RAG Accuracy (End-to-End QA)

 Evaluates how accurately the system extracts and interprets document content for downstream RAG tasks such as question answering.

#### **Key Metrics:**

- © Exact Match (EM) and F1-score on simple QA tasks.
- Interpolation Methods Building Buil

#### 📬 2. Processing Efficiency

 Measures how quickly and efficiently the system processes documents at scale.

#### **Key Metrics:**

- X Latency (seconds per document).
- Throughput (documents per second or pages per minute).
- Inference time for layout and extraction modules.





Reference: Intelligent Document Processing Leaderboard

idp-leaderboard.org

## SynGen AI – Synthetic Dataset Generation as a Service

Transforming Enterprise Documents into High-Fidelity Synthetic Instruction—Response Pairs for Small Language Model Training and RAG Optimization

#### S. S. Srinivas, Shivam Gupta, Akash Das

- StructiDoc AI transforms unstructured documents into structured data.
- Off-the-shelf SLMs are general-purpose and not fine-tuned on enterprise-specific tasks.
- SynGen AI provides high-fidelity synthetic data (QA pairs) from enterprise structured data tailored to customize SLMs for Domain-specific RAG over enterprise knowledge.

#### • 📥 Synthetic Dataset Creation Workflow:

- E Step 1: Input Source Collection
  - Collect enterprise documents and fed into StructiDoc AI platform, to obtain structured, machine-readable data optimized for AI workflows.
- 🗱 Step 2: Instruction-Response Pair Generation
  - Use large Vision-Language Models (e.g., GPT-4o, Gemini 2.5 Pro) to generate synthetic instruction-response pairs.
- Step 3: Dataset Typing for Specialized SLM Training
  - Verified Fact QA: Extracts fact-grounded answers from complex enterprise documents.
  - Reasoned Explanation QA: Generates step-by-step reasoning and explanations.
  - Contextual Grounding QA: Produces responses grounded in specific localized or distributed content segments.
- 4 Step 4: Quality Evaluation with Reward Models
  - Filter and validate responses using reward models like Nemotron-4-340B to ensure factuality, relevance, and completeness.
- Step 5: Student SLM Training via PORAG
  - Fine-tune smaller SLMs using Policy-Optimized RAG (PORAG) on these verified synthetic datasets for improved retrieval-augmented reasoning.

#### Outcome:

 Cost-effective, privacy-safe, and specialized training data powering high-fidelity RAG systems for enterprise document understanding.

# PORAG Trainer – Policy-Optimized Fine-Tuning as a Service

Leveraging Synthetic Instruction—Response Datasets to Fine-Tune Specialized SLMs for Accurate RAG

#### S. S. Srinivas, Shivam Gupta, Akash Das

- RL fine-tuning technique to customize SLMs on synthetic instruction–response datasets generated from enterprise data, optimizing the model to generate accurate responses for domain-specific RAG.
- It is essentially to teach the SLM to effectively utilize the retrieved context and generate responses grounded in enterprise knowledge.

- \$\mathbb{S}\$ 1. Inefficiency of Current Retrieval-Augmented Generation (RAG) Systems
  - RAG pipelines power many real-world applications (e.g., search, chatbots, copilots) by grounding LLM responses with external knowledge.
  - However, they suffer from:
    - ス Redundant or irrelevant retrievals, increasing latency (↑ ② ms/s per query).
    - Unnecessary computational overhead, lowering throughput (↓ tokens/s).
    - Key-Value (KV) caching overhead, increasing GPU VRAM usage (↑ GB). This limits the maximum context length and reduces batch size, increasing memory overhead and computational cost.
  - These limitations make RAG pipelines inefficient, costly, and difficult to scale in production environments.

#### Why Policy-Optimized Retrieval-Augmented Generation (PORAG) Is Relevant

- Problem: Existing RAG systems struggle with effective utilization of retrieved context
- Group Relative Policy Optimization (GRPO) is a reinforcement learning-based fine-tuning algorithm to enhance the reasoning capabilities of LLMs.
- Our Approach: Fine-tune SLMs through policy optimization over retrieved contexts
  - Integrates retrieval directly into the instruction tuning process
  - The "policy" refers to the SLM's parameters that govern text generation
  - Uses GRPO to update the language-only model parameters
  - Keeps retrieval mechanism fixed (computational efficiency)

#### The GRPO Loss Function:

- The GRPO loss function is a clipped policy optimization objective with group-relative advantage and KL penalty.
- Fine-tuning Efficiency: Uses QLoRA to reduce memory and compute overhead during training

#### Composite Reward Function for SLM Policy Optimization:

- Applies only to generated responses (retrieval is fixed)
- $R(y) = 0.3 \times \text{ROUGE-L F1} + 0.2 \times \text{Length Ratio Penalty} + 0.5 \times \text{LLM-as-Judge Score}$
- Optimizes for semantic similarity, brevity, factual correctness, and relevance

#### Key Benefits:

- Efficient Inference: Single-shot decoding with standard sampling
- No Multi-Candidate Ranking: Avoids expensive reward computation at inference
   Supplies Performance: Similificantly outpurforms usually PAC on fortunity matrix.
- Superior Performance: Significantly outperforms vanilla RAG on factuality metrics

## OptiInfer AI – Test-Time Inference Optimization as a Service

Optimizing Language Model Serving for Speed, Efficiency, and Scalability

S. S. Srinivas, Shivam Gupta, Akash Das, Venkataramana Runkana

- Optimizes language model serving for faster response generation without modifying model weights.
- Reduces latency, memory usage, and cost using system-level and reasoning-level optimizations.
- Enables scalable and efficient inference for production-grade RAG applications.

#### S Limitations of Static RAG:

- Unnecessary Retrievals: Always retrieves without checking if the available context is already sufficient, resulting in unnecessary latency and higher retrieval costs.
- Q Imprecise Querying: Builds a static query from the initial user input, missing
  opportunities to adapt queries based on evolving context or partial answers, leading
  to incomplete or inaccurate responses.
- Fixed Reasoning Depth: Uses static generation lengths, risking over-generation on simple tasks or under-generation on complex tasks.

#### ■ Adaptive Inference Optimization for RAG:

- It modifies the behavior at inference time by dynamically deciding when to retrieve and what to retrieve based on the evolving context during generation without altering the model weights.
- Context-Aware Querying: Leverages attention over the entire context to build precise, context-aware queries targeting missing information to fill information gaps to generate accurate response.
- Adaptive Reasoning: Varies generation depth based on task complexity and evolving context, balancing quality and efficiency.

#### 

Improves retrieval precision, reduces latency, and enhances factuality.

#### System-level Optimization Techniques

- We focus on low-level system-level optimizations that improve hardware-level performance to maximize runtime performance of SLMs.
- System-level optimizations focus on improving runtime efficiency of language models without modifying their parameters, targeting key system-level metrics: latency, throughput, and memory usage.

#### **Key Performance Metrics:**

- Latency: Time taken to generate a complete response (lower is better)
- Throughput: Number of tokens generated per second (higher is better)
- Memory Efficiency: GPU memory (VRAM) consumption impacting batch size and scalability

#### Techniques:

- FlashAttention: Efficient attention computation reduces memory bandwidth bottlenecks, improving both latency and throughput.
- PagedAttention with KV-Cache Quantization: Organizes the KV-cache into non-contiguous memory blocks to avoid fragmentation, improving memory efficiency and supporting larger batch sizes.
- Lookahead Decoding: Speculatively generates and verifies tokens in parallel to reduce generation latency while maintaining output quality.

#### Characteristics:

 Require no retraining or fine-tuning of model weights. Do not require multiple decoding passes, focus on accelerating vanilla decoding while maintaining output quality. Purely engineering/system-level improvements. 4日 → 4周 → 4 重 → 4 重 → 9 9 ○

#### $\ensuremath{{\mathbb Q}}$ Reasoning-level Optimization Techniques

At test time, algorithmic or reasoning-level optimizations can significantly improve the **factuality**, **reliability**, and **quality** of model outputs by modifying the generation strategy—without requiring any fine-tuning or retraining of model weights.

#### Key Focus Areas:

- Multi-Path Reasoning: Explore multiple reasoning trajectories and select the most consistent answer to improve robustness.
- Expert-Like Reflection: Simulate expert behaviors such as critique, reflection, and structured re-evaluation.

#### Core Characteristics:

- Works entirely at inference-time without modifying model weights.
- Focuses on improving output quality rather than computational speed.
- May increase computational cost by generating and evaluating multiple candidate responses.
- Relies on advanced decoding algorithms rather than parameter updates or retraining.

#### Benefits at a Glance:

- O Improve response quality without model fine-tuning.
- © Enhance factuality by verifying consistency across reasoning paths.
- Opnically control computational effort based on task complexity.
- Simulate expert-like critique and structured reasoning to improve reliability.

#### Advanced Inference Techniques

- Self-Consistency: Selects the most consistent answer by clustering multiple independently generated reasoning paths.
- Best-of-N Sampling: Picks the best from N candidates by self-evaluating response quality.
- Chain-of-Thought with Reflection: Guides reasoning through structured thinking, reflection, and answering phases in a single pass.
- Entropy-Guided Decoding: Dynamically adjusts sampling parameters based on model uncertainty to balance exploration and precision.
- Chain-of-Thought Decoding: Explores multiple reasoning paths and selects the most reliable based on token-level scoring.
- RE<sup>2</sup> (Re-Reading and Re-Analyzing): Structures reasoning into reading, re-reading, and final answer phases for deeper analysis.
- Mixture of Agents (MoA): Combines diverse generation, critique, and synthesis to produce refined responses.
- Reimplementation Then Optimize (RTO): Refines solutions by re-implementing from extracted specs and optimizing the final output.
- PlanSearch: Decomposes complex queries into multi-step planning and transformation stages before answering.
- Monte Carlo Tree Search (MCTS): Searches through reasoning paths using simulation and backpropagation for optimal responses.
- R\* Algorithm: Uses guided tree search with consistency checks to ensure reliable and structured reasoning. 4日 → 4周 → 4 重 → 4 重 → 9 9 ○

#### Policy-Optimized RAG (PORAG):

- RL fine-tuning technique for domain customization of SLMs
- Significantly improves factual accuracy in responses grounded in enterprise knowledge.

#### Adaptive Inference for RAG:

- Decides when to retrieve based on uncertainty or knowledge gaps during generation.
- Decides what to retrieve by building precise, context-aware queries targeting only missing information.
- Balances retrieval effort and response quality to reduce latency and cost without retraining the model.

#### Inference-Time Optimizations:

- Accelerates token generation and reduces memory usage
- Improves output quality through multi-path verification
- Enhances reliability through expert-like reasoning patterns
- Achieves better results without the need for fine-tuning

**Key Impact:** Enables faster, more accurate, and cost-efficient RAG across diverse applications and environments

#### % Unified AI Workflow: From Raw Enterprise Data to Optimized AI Outputs

- & Unlocking Value Across the Full AI Stack
  - Data Layer: Automate extraction, structuring, and synthetic data generation from enterprise documents.
  - Model Layer: Fine-tune SLMs with PORAG for domain-specific RAG optimization.
  - Inference Layer: Deploy scalable, cost-effective model serving with OptiInfer AI.

# Unstructured Docs StructiDoc Structured Data PDFs, Images JSON, Graphs Synthetic Data Inst-Resp Pairs





A Closed-Loop, Physics-Aware Agentic Framework for Auto-Generating Chemical Process

**and Instrumentation Diagrams** 

S. S. Srinivas, Shivam Gupta, Akash Das, Venkataramana Runkana

#### Process Flow Diagrams (PFDs):

- High-level schematic showing material and energy flows through production processes.
- Depict major equipment, process streams, and key operating conditions.
- Highlight what happens and where it happens in the process.

#### Piping and Instrumentation Diagrams (PIDs):

- Build upon PFDs by detailing valves, sensors, control loops, and actuators.
- Illustrate how the process operates and is controlled.
- Essential for safety, operational stability, and maintenance.

#### ● ▼ Foundation for Digital Twins and Al-Driven Automation

- PFDs and PIDs serve as foundational engineering schematics for digital twins.
- Digital twins integrate:
  - First-principles or data-driven methods.
  - Real-time sensor and actuator data streams from the physical process.
- Enable dynamic monitoring, predictive control, and closed-loop optimization.
- Power Al-driven automation across chemical manufacturing operations.

#### ■ Al Transforming Chemicals and Materials Discovery

- Generative AI accelerates discovery of:
  - Environmentally sustainable specialty chemicals.
  - Low-cost, high-performance materials-based products.
- Reduces dependency on expensive, slow lab-based trial-and-error workflows.
- Aids in simulation workflows and lowers R&D costs.
- Enables faster innovation and product development cycles.

#### ■ A Deployment Bottleneck: From Discovery to Production

- Challenge:Scaling discoveries from simulations or lab experiments requires developing new industrial production process.
- Current methods fail to:
  - Auto-generate industrial-scale PFDs and PIDs.
  - Justify design and control decisions.
  - Integrate physics-aware simulations for feasibility validation.
- Result: Slow and expertise-intensive manual workflows that limit scalability.

#### Ø Gaps in Process Context and Feasibility Validation

- Existing methods overlook:
  - High-level objectives and process sequencing.
  - Operational control, monitoring, and safety logic.
- Lack of simulator-based feasibility checks compromises reliability.
- These bottlenecks limits digital twin accuracy and scalable AI deployment.

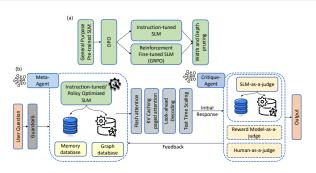
#### ■ © Closed-Loop, Self-Driving Lab Framework

- Cloud-based SaaS platform to automate:
  - High-fidelity PFD/PID generation.
  - Design, simulation, and optimization of process schematics with minimal human input.

#### Integrates Physics-Aware AI:

- Combines first-principles based process simulations with AI to ensure physical and operational feasibility.
- Validates generated schematics through reflection with simulator-based verification
- Provides continuous self-improvement through Al-driven feedback loops.
- Offers end-to-end process schematics modeling and validation.
- Expedites the simulation-to-lab-to-pilot-to-plant scale-up pipeline.
- Ensures that only viable, sustainable, and efficient processes advance to commercialization.

#### Overview of the Integrated Framework

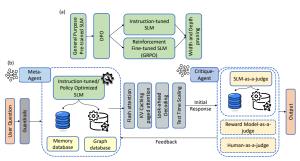


#### SLM Fine-Tuning Pipeline:

- Begins with Direct Preference Optimization (DPO) for alignment.
- Followed by Instruction Tuning or Group Relative Policy Optimization (GRPO).
- Concludes with optional width and depth pruning for efficiency.

#### Operational RAG Framework:

- Meta-Agent coordinates task planning and tool selection. Specialized SLM interacts with memory and graph databases for context retrieval.
- Inference accelerated using:
  - FlashAttention: Improves throughput and latency by reducing memory bandwidth bottlenecks in attention computation.



#### Core Optimizations:

- Paged KV Caching: Enhances memory efficiency by reducing memory fragmentation, enabling larger batch sizes during inference.
- Lookahead Decoding: Lowers latency by speculatively generating multiple tokens in parallel without sacrificing output quality.
- Test-Time Scaling: Increases factual accuracy by using techniques like multi-step reasoning, self-reflection, and re-ranking during inference.

#### • Iterative Response Refinement:

- Critique-Agent manages iterative feedback loops.
- Uses diverse judges:
  - Nemotron-4-340B reward model, SLM-as-a-judge, Human evaluations
- Offline-Third party Process simulations Verification.

- Agentic Multimodal AI for
- Hyper-Personalized B2B and B2C Advertising in Competitive Markets:
- An Al-Driven Competitive Advertising Framework

S. S. Srinivas, Shivam Gupta, Akash Das, Venkataramana Runkana

#### ▲ From Al-Led Product Innovation to Market Adoption

#### Industry Context:

- Al accelerates material discovery in energy, electronics, and FMCG.
- Success depends on bridging the gap between industrial scale-up and market adoption.

#### Commercialization Challenges:

- Weak value articulation limits market acceptance.
- New products struggle against established brands.
- One-size-fits-all messaging fails across regions.

#### Market Risks:

- Price wars from undifferentiated messaging.
- SKU-level cannibalization in product portfolios(Products from the same brand compete with each other, hurting overall sales).
- Low engagement from non-localized campaigns.

#### Limitations of Current Tools:

- Siloed R&D and marketing with no market feedback loop.
- Poor user modeling and static competitive insights.
- Lack of adaptive, privacy-safe campaign optimization.

#### Proposed Al-Driven Solution:

- Connect product-market fit with live market insights.
- Personalize engagement with data-driven targeting and adaptive creatives.
- Scale competitive messaging using agentic AI (MAAMS, PAG, CHPAS).

#### What is Programmatic Advertising?

- Automated buying and selling of digital ad space using software platforms.
- Replaces manual media buying with real-time, algorithmic auctions across digital channels.

#### • Advertising Channels in the Chemical Industry:

- Search: Google Ads, Bing Ads keyword-based auctions for buyer intent.
- Social: LinkedIn (B2B), Instagram, TikTok (B2C) audience-based ad delivery.
- E-commerce: Amazon, Alibaba, Knowde product discovery and lead generation.

#### Key Stakeholders:

- Publisher: Platform offering ad inventory (e.g., Google, Knowde).
- Advertiser: Chemical manufacturers promoting products to buyers.
- SSP (Supply-Side Platform): Manages and sells publisher inventory.
- DSP (Demand-Side Platform): Enables real-time bidding by advertisers.

#### How It Works:

- Powered by Real-Time Bidding (RTB)—an automated auction system for selecting the most relevant ad.
- RTB is triggered instantly when a user visits a digital property.



#### ② Real-Time Bidding (RTB) Across Channels

#### What is RTB?

- RTB is the auction engine behind programmatic advertising.
- It operates across search, social, and e-commerce channels.
- Each ad opportunity—called an impression—is evaluated and sold in real time, typically within 100 milliseconds.

#### How RTB Works:

- A user visits a digital property (e.g., Google Search, LinkedIn feed, Knowde product page).
- The SSP offers the impression to multiple DSPs.
- OSPs evaluate:
  - User signals (location, industry, behavior).
  - Page context (chemical category, product detail).
  - Channel data (search intent, social engagement).
- 4 Advertisers submit real-time bids.
- The highest bidder wins; their ad is instantly shown.

#### • Auction Types:

- Header Bidding: A pre-auction strategy where multiple SSPs bid simultaneously for the same impression—maximizing competition and publisher revenue.
- Second-Price Auction: A pricing mechanism where the highest bidder wins, but pays only the second-highest bid—ensuring fairness and cost-efficiency for advertisers.

#### ♣ Auction Mechanics in Programmatic Advertising

#### Header Bidding (Pre-Auction Strategy):

- A pre-auction mechanism where multiple SSPs bid simultaneously.
- Replaces the "waterfall" approach, where only one SSP (e.g., Google Ad Manager) gets priority.
- The highest SSP bid is passed to the final ad server auction.

#### Why it matters:

- Increases competition across SSPs (e.g., Google, Amazon, Xandr).
- Boosts publisher revenue and avoids walled garden dominance.

#### First-Price Auction:

- The winner pays exactly what they bid.
- Leads to strategic bidding, not truthful bidding—creates risk and inefficiency.
- Second-Price Auction (used in RTB):
  - The highest bidder wins, but pays only the second-highest bid + \$0.01.
  - Encourages truthful bidding—advertisers can bid their actual maximum value.
  - Reduces overpayment risk, improving auction fairness and pricing efficiency.

#### Data-Driven Targeting for Smarter Bidding

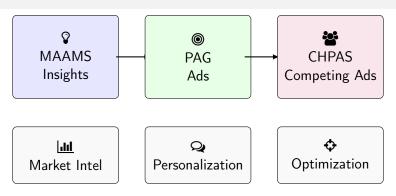
- Real-Time Bidding (RTB) selects ads in milliseconds, but effectiveness depends on targeting the right users.
- Ads shown to low-intent users waste budget—even if they win the auction.
- Precise targeting turns a fast auction into a smart investment by improving return on ad spend (ROAS).
- Data-driven targeting uses first-party, third-party, and contextual signals to identify high-intent users.
- Targeting Inputs:
  - First-Party Data: From CRM systems and website behavior.
    - Customer profiles, purchase history, product preferences.
    - Website actions like product views, downloads, sample requests.
  - Third-Party Data: From Data Management Platforms (DMPs).
    - Industry-specific segments (e.g., plant managers, chemical buyers).
    - Facility type, production capability, and procurement roles.
  - Contextual and Geographic Signals:
    - Page relevance (e.g., chemical applications, compliance).
    - Location-based filters (e.g., logistics feasibility).
    - Privacy-aware processing (GDPR, CCPA compliant).
- These inputs enable DSPs to bid intelligently—based not just on price, but on the likelihood of conversion—across search, social, and e-commerce channels.

#### Personalizing the Ad After Winning the Bid

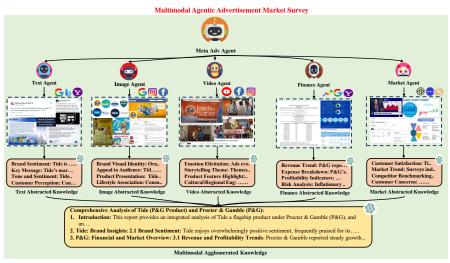
- Winning the bid secures ad space—but performance depends on delivering a relevant and personalized creative to the user.
- MI -driven creative selection considers:
  - User attributes (profile, behavior, intent).
  - Channel context (search, social, e-commerce).
  - Historical performance (CTR, conversions).
- Format varies by channel:
  - Static banners in search.
  - Videos or carousels in social feeds
- Dynamic Creative Optimization (DCO) personalizes:
  - Product visuals e.g., lab equipment for scientists vs. bulk packaging for procurement teams.
  - Messaging tone e.g., "improve yield and purity" for R&D vs. "reduce cost and downtime" for operations.
  - Highlighted benefits e.g., regulatory compliance for pharma vs. sustainability metrics for specialty chemicals.
- Performance metrics (clicks, downloads, sample requests) feed back into future targeting and creative optimization.
- Takeaway: Personalizing the ad creative after winning the bid is essential—the right message, delivered to the right user, drives outcomes. Winning the bid alone is not enough.

- Limitations of Traditional Approaches:
  - Lack of dynamic creative optimization to personalize engagement at scale.
  - Weak modeling of user personas and multimodal buyer behavior.
  - Inability to differentiate similar products from competing brands or within the same portfolio.
  - Poor adaptation to privacy regulations like GDPR and CCPA in real-time targeting workflows.
- Coptimization Goals for Ad Performance:
  - Maximize Return on Ad Spend (ROAS) by optimizing targeting, messaging, and acquisition efficiency.
  - SKU (Stock Keeping Unit) cannibalization occurs when multiple products from the same company compete for the same customer, reducing total sales instead of expanding market share.
  - Ensure privacy-compliant optimization using real-time market signals.
- ♥ Why Advertising Matters in Chemical GTM(Go-To-Market):
  - Moves products beyond lab-scale innovation to real market adoption.
  - Clarifies unique technical and commercial value to non-technical buyers.
  - Helps new launches stand out from entrenched incumbents.
  - Reaches niche B2B and broad B2C audiences across global markets.
  - Drives measurable outcomes through adaptive, data-driven campaigns.

#### Agentic Al Advertisement Framework

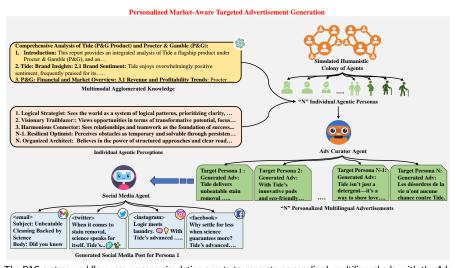


- MAAMS (Multimodal Agentic Advertisement Market Survey): Surveys the market using specialized agents to gather multimodal intelligence on brand perception, financials, and competitor positioning.
- PAG (Personalized Market-Aware Targeted Advertisement Generation): Generates personalized, multilingual ads by simulating diverse consumer personas and tailoring content to their specific preferences and cultural contexts.
- CHPAS (Competitive Hyper-Personalized Advertisement System): Optimizes these
  personalized ads for competitive scenarios by strategically highlighting unique selling
  points against rival products to maximize relevance and effectiveness for each target user.



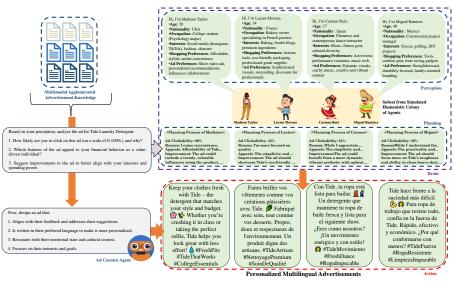
The MAAMS system aggregates multimodal insights via a Meta-Agent to analyze brand sentiment, performance, and market standing.

#### Personalized Market-Aware Targeted Advertisement Generation(PAG)



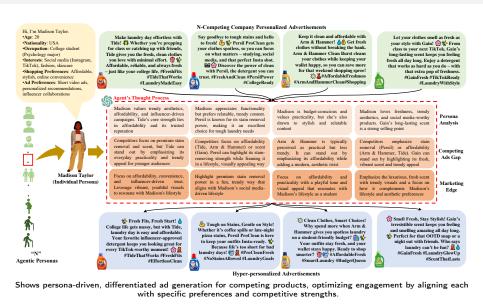
The PAG system workflow uses persona-simulating agents to generate personalized, multilingual ads, with the Adv Curator ensuring cultural fit and the Social Media Agent optimizing platform delivery.

## Personalized & Multilingual Ad Creation



The PAG system generates personalized, multilingual ads by evaluating consumer preferences, cultural context, and feedback using multimodal knowledge to maximize engagement and drive action.

#### Hyper-Personalized Competitive Ads



イロト (個) (4) (4) (4) (4) (4) (4)

## Ad Copy Optimization for Consumer Chemical Products

- Objective:
  - Enhance the effectiveness of Al-generated ad copy for consumer chemical products (e.g., personal care, home care, nutrition, and cleaning).
  - Maximize emotional impact, brand recall, and purchase intent across diverse consumer segments.
- Evaluation Framework:
  - Based on CHPAS: Competitive Hyper-Personalized Advertisement System.
  - Uses Simulated Humanistic Agentic Personas (SHAP) to assess ad variations across key consumer dimensions:
    - Emotional Impact Evaluator Assesses warmth, empathy, excitement, or trust.
    - Persona Resonance Filter Measures alignment with different lifestyle profiles (e.g., eco-conscious parent, busy professional, wellness seeker).
    - Q Cultural Context Checker Ensures relevance across regions, languages, and norms
    - Brand-Value Alignment Validates messaging against brand positioning (e.g., sustainability, safety, luxury).
  - Scoring Dimensions (0-10 scale): Emotional Resonance, Brand Fit, Clarity, Cultural Relevance
- Experimental Setup:
  - Run synthetic evaluations across channels like Instagram, TikTok, Meta, and e-commerce banners.
  - Test copy styles: benefit-led vs. feature-led, sensory words, visual callouts, and CTA phrasing. = 900 €

Motivation: Real-world A/B testing for ad personalization is costly, risky, and constrained by privacy laws (e.g., GDPR, CCPA).

Goal: Develop a privacy-compliant framework to simulate and optimize ad strategies across competing products—entirely offline.

#### Pipeline Overview:

- Persona Profiling: Generate synthetic user profiles with demographic and behavioral traits.
- ② ♥ Product Modeling: Evaluate product features, pricing, and brand perception.
- ⑤ ☼ Competitive Simulation: Model competing products and market scenarios.
- Product-Persona Alignment: Score alignment between personas and product attributes.
- **⑤** ★ Ad Generation: Create generic and persona-optimized ads using LLMs.
- § LLM Evaluation: Score ads for relevance, persuasion, and emotional impact (e.g., GPT-4 Omni, Nemotron).

Outcome: A scalable, cost-effective framework enabling hyper-personalized ad optimization without real-world deployment.

# ChemConnect Al:

Digitizing Chemical Commerce with Agentic Multimodal B2B Marketplace Intelligence

📜 Unified Marketplace | 🗱 Al Product Structuring | 🔼 Data-Driven Competitive Insights

S. S. Srinivas, Shivam Gupta, Akash Das, Venkataramana Runkana

July 4, 2025

#### Why It Is Needed

- The \$20T+ global chemical industry remains highly fragmented, with <5% of transactions conducted digitally.
- ■ Relies on legacy mechanisms—emails, trade shows, and phone calls—that slow
   procurement and discovery cycles.
- Supplier data scattered across unstructured PDFs, spreadsheets, and siloed systems.
- Q No centralized platform for technical product discovery, structured comparison. sampling, and quoting

#### Why It Matters

- 1 Improves sourcing efficiency, reducing time-to-market and costs.
- Ecentralizes supplier catalogs into machine-readable structured data for search and discovery.
- Pigitally connects verified suppliers to high-intent global buyers.
- Accelerates R&D by simplifying material discovery, sampling, and quoting.

#### ⚠ What Is ChemConnect AI?

- A global B2B digital marketplace for chemicals, polymers, ingredients and etc.
- Leverages AI to extract, normalize, and structure unstandardized supplier data for technical searchability.
- Provides seamless search, sampling, quoting, and procurement workflows.
  - In B2B chemical marketplaces, sampling refers to the process where buyers request free or paid small quantities of materials (called samples) to evaluate them before committing to bulk purchases.
  - It helps R&D teams, formulators, and procurement professionals to:
    - Test performance in their own formulations or manufacturing processes.
    - Validate quality and specifications against their requirements.
    - Compare multiple suppliers' offerings before finalizing large-scale procurement.
- Powers digital transformation across the chemical industry by connecting suppliers and buyers on a unified platform.
- Empowers suppliers to go digital through branded storefronts powered by ChemConnect's backend

#### What ChemConnect Al Offers

- \(\begin{align\*}\) Marketplace: Centralized, searchable catalog of verified suppliers and products.
- Al Engine: Extracts product specs from unstructured data (PDFs, Excel).
- Master Data Management (MDM): Ensures cross-supplier consistency and normalization
- Qustomer Experience Platform (CXP): Powers branded storefronts with Al-driven search, filtering, and sample/quote workflows.
- ChemConnect AI offers more than just hosting products on its marketplace.
- It enables suppliers to launch their own branded digital storefronts on their company websites
- These storefronts use ChemConnect AI's technology to provide:
  - Advanced product search and filtering.
  - Structured data presentation for easy discovery.
  - Lead capture through sample and quote requests.
  - Access to product documents, certifications, and specifications.
- This helps suppliers digitally engage customers directly on their own websites while leveraging ChemConnect AI's platform capabilities.

#### Stakeholders

- Suppliers: Upload and structure chemical products with specs, documents, and certifications.
- Buyers: R&D labs, formulators, and procurement teams seeking discovery, sampling, and technical procurement.

#### ∆i∆ Mutual Value

- Suppliers: Expand reach, digitize catalogs, launch storefronts, and convert leads to deals.
- Buyers: Perform side-by-side comparisons, access verified data, and initiate sampling or guotes with fewer intermediaries.

## Why Not Amazon, Flipkart, or Uber?

- Consumer E-commerce: Lacks multi-attribute, compliance-focused technical filtering.
- Service Marketplaces: Not designed for product procurement, sampling, or data compliance.
- ChemConnect AI: Purpose-built for industrial-grade search, domain-specific ontology, and technical workflows

# ChemConnect Insights Platform

### ■ Introducing ChemConnect Insights:

- Subscription-based analytics platform for chemical suppliers and industrial buyers.
- Provides real-time, sector-specific, and behavioral market intelligence to drive data-informed decisions.

## ■ Sector & Regional Market Intelligence:

- Access industry-specific demand insights across coatings, pharma, and etc.
  - Enable targeted go-to-market, pricing, and distribution strategies across global regions and sectors.

# ② Real-Time Engagement Analytics:

- Track product views, sampling requests, quote submissions, and buyer engagement.
- It allows suppliers to optimize product positioning, catalog performance, and customer outreach dynamically.

# • Page ChemConnect GPT (Al-Powered B2B Insights Assistant):

- Al assistant( chat-based tool) trained on ChemConnect marketplace data helps suppliers and buyers answer performance queries, generate competitive benchmarks.
- Provides actionable sales or sourcing recommendations.

#### Suppliers:

- Understand how their products are performing (e.g., view rates, quote requests, conversion rates).
- Compare their performance to similar competitors (competitive benchmarks).

#### Buyers

- Discover better or cheaper alternatives based on historical sourcing data.
- Get recommendations on which suppliers or products fit their needs based on marketplace insights. **■** 990

#### III Behavioral Metrics Dashboard:

- Monitor buyer retention, repeat sampling behaviors, and supplier share-of-voice in in platform search and category placements.
  - Buyer Retention: How many buyers come back to the same supplier or product over time.
  - Repeat Sampling: How often the same buyers request new samples of the same product or related products.
  - Share-of-Voice: How prominently a supplier's products appear in search results and category listings compared to competitors.
- It offers a comprehensive view of customer engagement across the marketplace.

### \$ Strategic Business Impact:

- Provides data-backed decision support for suppliers and buyers.
- Unlocks a recurring SaaS revenue stream alongside marketplace transactions.
  - (i) charge a commission or fee on every successful B2B sale made through your digital marketplace.
  - (ii) Suppliers or buyers pay monthly or yearly to access advanced analytics).



Questions & Discussion