

A Comprehensive Analysis of Linear Algebra-Based Performance Modeling and Enterprise Invoice Processing

Goutam Parashuram Gotur¹, Dr. E. Saravana Kumar²

¹ Department of Computer Science and Engineering, The Oxford College of Engineering, Bengaluru, Karnataka, India

² Department of Computer Science and Engineering, The Oxford College of Engineering, Bengaluru, Karnataka, India

Correspondence: goutamgotur2006@outlook.com, saraninfo@gmail.com

Abstract

Hybrid artificial intelligence architectures combining traditional computational methods with neural network residuals represent a paradigm shift in addressing complex real-world challenges. This report synthesizes two complementary approaches: (1) linear algebra-based digital system performance modeling that leverages matrix-vector operations enhanced with neural network approximators, and (2) optical character recognition (OCR) integrated with large language models (LLMs) for automated invoice processing in e-commerce environments. Both methodologies exemplify the principle of interpretability-efficiency trade-offs in modern AI systems. This work demonstrates how decomposing complex problems into interpretable baselines with neural residuals yields superior performance in accuracy, inference speed, and scalability compared to monolithic deep learning approaches. The report presents mathematical formulations, implementation strategies, empirical validation, and practical deployment considerations across diverse application domains.

Keywords: Hybrid AI systems, Linear Algebra, Neural Networks, Large Language Models, OCR, Performance Modeling, Automated data extraction, E-commerce, Interpretability, Scalability

1 Introduction

Digital systems and enterprise automation present increasingly complex challenges that resist simple linear solutions. Traditional approaches—whether purely mathematical (linear algebra-based) or purely data-driven (deep neural networks)—have inherent limitations. Linear models Ax , while computationally efficient and interpretable, frequently underfit nonlinear phenomena such as cache effects, queueing dynamics, and contention in hardware systems. Conversely, monolithic neural networks achieve high accuracy but at the cost of computational expense, reduced interpretability, and greater sample complexity.

The emergence of hybrid architectures addresses these limitations by combining the strengths of each paradigm. In the domain of digital system performance prediction, Gotur and Shree propose decomposing the true performance mapping $f(x)$ into an interpretable linear baseline Ax and a compact neural residual $\hat{f}(x; \theta)$, formulated as:

$$y(x) = Ax + \hat{f}(x; \theta)$$

This decomposition enables practitioners to maintain computational efficiency and interpretability while leveraging neural networks to capture residual nonlinearities .

Similarly, in enterprise invoice processing, Patel and Pandit address the challenge of diverse document formats—PDFs, handwritten documents, and scanned JPG images—by combining optical character recognition (OCR) technology with generative AI models. Rather than attempting direct end-to-end learning from raw images to extracted entities, their approach decomposes the problem: OCR converts image-based invoices into machine-readable text, and a large language model then performs entity extraction and standardization .

Both methodologies exemplify a broader principle: complex real-world problems benefit from hierarchical decomposition, where each layer handles phenomena at appropriate levels of abstraction. This report synthesizes these two case studies to articulate common principles, mathematical foundations, implementation strategies, and insights for practitioners developing hybrid AI systems.

1.1 Problem Context and Motivation

Digital System Performance Prediction: Organizations must predict system performance (latency, throughput, tail metrics) to optimize resource allocation, diagnose bottlenecks, and ensure service level agreements . A purely linear model may capture first-order resource contributions but miss nonlinear phenomena. A large neural network may achieve high accuracy but become a black box, complicating diagnosis and requiring extensive retraining when system architectures change .

Enterprise Invoice Processing: E-commerce platforms process invoices from thousands of vendors, each submitting documents in different formats and layouts . Manual extraction of key entities (invoice number, vendor name, total amount, tax, line items) does not scale. Purely automated pipelines using naive text extraction fail on handwritten and scanned documents. An integrated solution combining OCR and LLMs reduces manual intervention and improves accuracy .

1.2 Contributions and Report Structure

This report makes the following contributions:

1. **Unified Framework:** Articulates a common design principle—interpretable baseline plus neural residual—applicable across multiple domains.
2. **Mathematical Formulation:** Formalizes both approaches with rigorous notation, error bounds, and complexity analysis.
3. **Implementation Guidance:** Provides practical algorithms, hyperparameter recommendations, and deployment strategies.
4. **Empirical Evidence:** Synthesizes results from both case studies demonstrating improvements in accuracy, inference speed, and scalability.
5. **Future Directions:** Identifies limitations and proposes extensions, including transformer-based residuals and uncertainty quantification.

The report proceeds as follows: Section 2 presents the mathematical foundations of hybrid linear-neural models for performance prediction. Section 3 describes invoice processing via OCR and LLM integration. Section 4 discusses common principles and design trade-offs. Section 5 addresses implementation and deployment. Section 6 concludes with limitations and future work.

2 Hybrid Linear-Neural Models for Digital System Performance

2.1 Mathematical Formulation

Following the framework developed by Goutam , we formalize the hybrid performance prediction model.

Definition 1 (Hybrid Performance Model): Let $x \in \mathbb{R}^n$ be a workload feature vector encoding arrival rates, request sizes, concurrency levels, and other characteristics of the workload. Let $y \in \mathbb{R}^m$ be the observed performance vector, typically containing latencies, throughput metrics, and tail latency percentiles. The hybrid model is defined as:

$$y(x) = Ax + \hat{f}(x; \theta)$$

where:

- $A \in \mathbb{R}^{m \times n}$ is a learnable linear baseline matrix capturing per-resource contributions
- $\hat{f}(x; \theta)$ is a neural network with parameters θ that approximate residual nonlinearities
- $f(x)$ represents the true (unknown) performance mapping

Baseline Model: The linear component alone is:

$$y_0(x) = Ax$$

Residual Definition: The residual function is:

$$r(x) := f(x) - Ax$$

The neural network aims to approximate this residual: $\hat{f}(x; \theta) \approx r(x)$.

2.2 Neural Network Architecture

The neural residual is implemented as a standard feedforward network with L layers :

$$h_1 = \sigma_1(W_1x + b_1)$$

$$h_\ell = \sigma_\ell(W_\ell h_{\ell-1} + b_\ell), \ell = 2, \dots, L - 1$$

$$\hat{f}(x; \theta) = W_L h_{L-1} + b_L$$

where σ_ℓ are activation functions (ReLU, GELU, etc.), and $\theta = \{W_\ell, b_\ell\}_{\ell=1}^L$ collects all parameters.

Architecture Rationale: Typical configurations use 1–3 hidden layers with widths of 32–256 neurons . The residual network should be compact, just large enough to capture nonlinearities left unexplained by the linear baseline. This contrasts with monolithic neural networks, which require much larger capacity (e.g., 200k parameters), as shown in Table I .

2.3 Training Objective

Joint optimization of A and neural parameters θ is formulated as :

$$\min_{A, \theta} \frac{1}{N} \sum_{i=1}^N \ell(y^{(i)}, Ax^{(i)} + \hat{f}(x^{(i)}; \theta)) + \lambda(\|A\|_F^2 + \Omega(\theta))$$

where:

- ℓ is a loss function (e.g., mean squared error for regression)
- $\Omega(\theta)$ is a regularization term (weight decay, dropout rate, etc.)
- λ is the overall regularization coefficient
- N is the number of training samples

2.4 Error Analysis

Theorem 1 (Error Decomposition): The prediction error of the hybrid model is bounded as:

$$\|E(x)\| \leq \|f(x) - Ax\| + \|\hat{f}(x; \theta) - (f(x) - Ax)\|$$

where $E(x) := f(x) - (Ax + \hat{f}(x; \theta))$ is the total prediction error .

Interpretation: The total error decomposes into two components:

1. **Baseline Bias:** $\|f(x) - Ax\|$ —the error introduced by the linear approximation
2. **Neural Approximation Error:** $\|\hat{f}(x; \theta) - r(x)\|$ —the residual network's approximation error

This decomposition provides interpretability: practitioners can analyze which phenomena are explained by linear baseline and which require neural correction.

Corollary (Error Bound on Compact Domain): If \hat{f} approximates the residual $r = f - Ax$ within ϵ on a compact domain \mathcal{X} , then:

$$\sup_{x \in \mathcal{X}} \|E(x)\| \leq \sup_{x \in \mathcal{X}} \|f(x) - Ax\| + \epsilon$$

By the universal approximation property of neural networks, provided sufficient capacity, ϵ can be made arbitrarily small.

2.5 Computational Complexity

Theorem 2 (Hybrid Complexity): The computational cost of hybrid inference is:

$$C_{\text{hybrid}} = \mathcal{O}(\text{nnz}(A)) + \sum_{\ell=1}^L d_{\ell-1}d_{\ell}$$

where:

- $\text{nnz}(A)$ is the number of nonzeros in matrix A
- $d_0 = n, d_1, \dots, d_L = m$ are layer dimensions
- The first term is matrix-vector multiplication cost
- The second term is neural network forward pass cost

Practical Implications: This decomposition enables deployment trade-offs. Maintaining a sparse A (few nonzeros) and small network widths reduces overall inference cost significantly compared to a monolithic neural network. For example, Gotur and Shree report :

- Linear-only model: 0.12 ms inference time, 1.2k parameters, MSE = 0.045
- Large neural network only: 1.8 ms inference time, 200k parameters, MSE = 0.012
- Hybrid model (proposed): 0.18 ms inference time, 8.5k parameters, MSE = 0.014

The hybrid approach achieves accuracy nearly as good as the large NN while maintaining inference speed close to linear-only and parameter count $24\times$ smaller than the monolithic NN .

2.6 Construction of the Linear Baseline

Several strategies exist for constructing the baseline matrix A :

Strategy 1: Ridge Initialization

$$A_{\text{ridge}} = YX^T (XX^T + \alpha I)^{-1}$$

where $Y \in \mathbb{R}^{m \times N}$ and $X \in \mathbb{R}^{n \times N}$ contain observed performance and features, and α is a regularization parameter.

Strategy 2: Sparse Baseline

$$\min_A \|Y - AX\|_F^2 + \lambda_1 \|A\|_1$$

promotes sparsity, making A more interpretable and computationally efficient.

Strategy 3: Low-Rank Factorization

$$A = UV^T, U \in \mathbb{R}^{m \times k}, V \in \mathbb{R}^{n \times k}$$

reduces the number of parameters, especially useful when $\min(m, n)$ is large.

2.7 Application Domain: CPU Scheduling and Memory Bandwidth

Gotur and Shree illustrate their approach with two case studies :

CPU Scheduling: Features include arrival rates, burst times, and priority levels. The linear component A models scheduling delays as linear functions of these features. The neural residual \hat{f} captures cache effects and context-switch overhead—inherently nonlinear phenomena influenced by system state.

Memory Bandwidth Prediction: Features encode read/write ratio, stride patterns, and concurrency. The baseline A provides a coarse bandwidth estimate; the neural residual models queuing for delays, contention effects, and prefetch optimization—all nonlinear and state-dependent.

3 Large Language Models and OCR for Automated Invoice Processing

3.1 Problem Context and Challenges

E-commerce platforms process invoices from thousands of vendors globally. Invoice formats vary widely:

- **PDF:** Digital invoices with complex layouts, multiple columns, embedded images, and variable structure
- **Handwritten:** Manual invoices with legibility and consistency challenges
- **Scanned/Image:** JPG snapshots of physical invoices, subject to lighting and angle variations

Key Challenges :

1. **Format Heterogeneity:** No standardized invoice template; layouts differ across vendors
2. **Manual Processing Bottleneck:** Employees manually extracting data from handwritten and scanned invoices, leading to errors and inefficiency
3. **OCR Limitations:** Handwritten invoices present variable handwriting styles, reducing OCR accuracy
4. **Layout Inconsistency:** Even structured PDFs vary in data arrangement, requiring flexible extraction logic
5. **Scalability:** Manual processing does not scale as vendor base grows

Operational Impact:

- High error rates in data extraction
- Reduced operational efficiency
- Inability to scale to tens of thousands of daily invoices
- Compliance and audit risks due to inaccuracies

3.2 Hybrid Solution: OCR + LLM Integration

Patel and Pandit propose a two-stage architecture decomposing the problem:

Invoice Processing Pipeline: OCR + LLM Integration

\bigskip

Stage 1: Input Formats (PDF, JPG, Handwritten)

↓

Stage 2: Optical Character Recognition (OCR)

↓

Stage 3: Text Preprocessing and Metadata Reading

↓

Stage 4: Feature Extraction and Standardization

↓

Stage 5: GenAI Entity Extraction Model (LLM)

↓

Stage 6: Storage in Vector Database

↓

Stage 7: Human Validation (Quality Assurance)

↓

Stage 8: Feedback Loop to Improve Model

↓

Stage 9: Structured Output for Business Applications

Figure 1: End-to-end architecture for centralized invoice processing combining OCR and LLM technologies

Stage 1: Input Formats

The system accepts invoices in PDF, scanned image (JPG), and document formats, each with different structural and semantic characteristics .

Stage 2: Optical Character Recognition

OCR converts handwritten and image-based invoices into machine-readable text, enabling downstream digital processing . OCR is particularly critical for:

- Handwritten invoices with variable legibility
- Scanned documents degraded by lighting and angle
- Legacy documents not digitally born

Stage 3: Text Preprocessing

After OCR, the system reads invoice metadata (date, vendor, payment terms) and ensures accurate feature extraction .

Stage 4: Feature and Property Extraction

The system identifies and extracts critical features:

- Invoice number and date
- Vendor details (name, account number, contact information)
- Tax amounts and applicable tax codes
- Line items (product descriptions, quantities, unit prices, extended amounts)
- Payment terms and due dates

- Total amount due

Stage 5: GenAI Entity Extraction Model

A large language model processes the digitized text to automatically extract and standardize entities . The LLM is trained in diverse invoice layouts, enabling robust generalization across varied formats. Key advantages:

- Handles varied layouts and structures
- Recognizes semantic relationships (e.g., which vendor information applies to which invoice)
- Performs implicit disambiguation and correction
- Scalable to new vendor types without manual rule engineering

Stage 6: Vector Database Storage

Extracted entities are stored in a vector database, enabling rapid semantic search, auditing, and downstream analytics .

Stage 7: Human Validation

Human reviewers perform quality assurance, checking for:

- Extraction accuracy
- Handling of ambiguous or unclear entries
- Edge cases not well-handled by automated systems

Stage 8: Feedback Loop

Human validation results are fed back to the LLM model, enabling continuous improvement . This active learning approach ensures the model adapts to new patterns and improves generalization.

Stage 9: Structured Output

The validated, structured data is output for use in:

- Financial reporting and reconciliation
- Tax compliance
- Audit trails
- Analytics and vendor performance tracking

3.3 Entity Extraction with LLM

The core of the invoice processing system is the LLM-based entity extraction model. Unlike traditional rule-based systems or simpler pattern-matching approaches, the LLM leverages:

1. **Semantic Understanding:** Contextual comprehension of invoice content, not just pattern matching
2. **Robustness to Layout Variation:** Ability to extract information from diverse document structures
3. **Generalization:** Training on diverse invoices enables application to unseen vendor formats

4. **Interpretability:** LLM can provide confidence scores and reasoning for extractions

Formulation: Given OCR-extracted text t from an invoice (PDF or handwritten), the LLM computers:

$$E = \text{LLM}_{\text{invoice}}(t; W)$$

where:

- E is a structured entity dictionary containing extracted fields
- W represents LLM weights, typically a large foundation model fine-tuned on invoice data
- The LLM may be prompted to output structured JSON, YAML, or tabular formats for downstream processing

3.4 Empirical Results and Impact

Patel and Pandit report significant operational improvements :

Efficiency Gains:

- 90% reduction in manual effort for invoice processing
- Automated processing of tens of thousands of invoices daily
- Elimination of per-invoice labor costs at scale

Accuracy Improvements:

- 75% reduction in data extraction errors compared to manual processing
- High accuracy across different invoice formats (PDFs, handwritten, scanned)
- Particularly effective for handwritten invoices (70% accuracy improvement vs. manual)

Scalability:

- No additional labor costs as platform grows
- Ability to process 100,000+ invoices daily
- Extensible to new vendor types without manual recoding

Centralized Data Management:

- All extracted data stored in standardized format
- Streamlined reporting and financial audits
- Regulatory compliance (e.g., tax reporting)

Use Case: A major e-commerce platform handling 200,000 vendors worldwide achieves :

- Processing of 100,000 invoices per day
- 70% improvement in handwritten invoice accuracy

- Drastic reduction in manual data entry errors
- Improved compliance and auditability

4 Unified Principles and Design Trade-Offs

4.1 Hybrid Architecture Philosophy

Both case studies exemplify a common architectural principle: **hierarchical decomposition with interpretable baselines and learned residuals**. This principal manifests differently in each domain but reveals deep insights applicable broadly.

Common Pattern :

1. **Problem Decomposition:** Break the complex mapping (performance prediction or invoice processing) into interpretable components plus a learned residual
2. **Baseline Implementation:** Encode domain knowledge as an interpretable, computationally cheap baseline
3. **Residual Learning:** Use flexible (neural) components to capture phenomena not captured by the baseline
4. **Integration:** Combine both components additively or multiplicatively, depending on problem structure

Advantages of this approach :

- **Interpretability:** Baseline remains transparent; practitioners understand what each component contributes
- **Sample Efficiency:** Residual learning is easier than learning the full mapping; fewer data required
- **Computational Efficiency:** Baseline is cheap; learned component can be compact
- **Modularity:** Each component can be updated, replaced, or analyzed independently
- **Debugging:** When predictions are wrong, practitioners can diagnose whether the error is in the baseline or the residual

4.2 Interpretability vs. Accuracy Trade-Off

There exists a fundamental trade-off between model interpretability and accuracy :

- **Purely Linear/Rule-Based:** Highly interpretable, computationally cheap, but often underfits (e.g., MSE = 0.045 for linear-only performance model)
- **Purely Neural/Deep:** High accuracy (e.g., MSE = 0.012 for large NN), but black-box nature and high computational cost

- **Hybrid:** Intermediate accuracy (MSE = 0.014, nearly as good as pure NN) with interpretability maintained via the transparent baseline

Hybrid models represent a principled middle ground, enabling practitioners to achieve near-state-of-the-art accuracy while preserving interpretability and computational efficiency.

4.3 Scalability Considerations

Digital System Performance: The hybrid approach enables cost-efficient scaling :

- Inference cost decomposes as $\mathcal{O}(\text{nnz}(A)) + \text{NN cost}$
- By maintaining sparse A and compact neural networks, overall cost scales sublinearly with model capacity

Invoice Processing: The OCR + LLM approach enables operational scaling:

- OCR preprocessing can be parallelized across documents
- LLM inference on OCR output is more efficient than end-to-end image-to-entity extraction
- Human validation can be sampled (e.g., 1–5% of invoices) to detect systematic errors while maintaining scalability

4.4 Generalization and Transfer Learning

Performance Modeling: Different systems (CPUs, memory systems, databases) have different nonlinear characteristics. The hybrid approach facilitates transfer learning:

- The baseline A can be retrained quickly on a new system
- The residual network may generalize across systems if nonlinearities are similar
- Alternatively, the baseline can transfer, and only the residual needs retraining

Invoice Processing: Different vendors submit invoices in different formats. The hybrid approach enables generalization :

- OCR is vendor-agnostic; the same OCR works across handwriting styles and scanned qualities
- LLM fine-tuning on a representative sample of vendor invoices enables generalization to new vendors
- Feedback loop continuously improves generalization

5 Implementation and Deployment

5.1 Data Preparation and Preprocessing

5.1.1 Performance Modeling Preprocessing

For digital system performance prediction, preprocessing includes :

Standardization: Normalize input features to zero mean and unit variance:

$$x' = \frac{x - \mu}{\sigma}$$

where μ and σ are computed on training data.

Feature Engineering: Domain-driven feature selection and derivation improves baseline linearity :

- Arrival rates, request sizes, burst durations
- Read/write ratios, stride patterns, cache behavior
- Concurrency levels, queue depths
- Hardware resource utilization percentages

Train-Validation-Test Split: Use 70% training, 15% validation, 15% test data .

5.1.2 Invoice Processing Data Preparation

For invoice systems, preprocessing includes :

Format Standardization: Convert all input documents (PDFs, images, handwritten scans) to a common format (e.g., page images) before OCR.

OCR Preprocessing:

- Image enhancement (brightness, contrast correction)
- Rotation normalization
- Denoising to improve legibility

Text Cleaning:

- Removing spurious line breaks
- Standardizing whitespace
- Preserving document structure (tables, lists)

Annotation: A sample of invoices (e.g., 200–500) should be manually annotated with ground truth entities for model training and evaluation.

5.2 Training Strategy for Hybrid Models

Algorithm 1: Staged Training for Linear-Neural Hybrid Models

Input: Training data $X \in \mathbb{R}^{n \times N}$, $Y \in \mathbb{R}^{m \times N}$

Hyperparameters: λ , learning rates η_A , η_θ , epochs

1. Initialize A via ridge regression:
 $A \leftarrow Y X^T (X X^T + \lambda I)^{-1}$

2. Stage 1 - Residual Learning:
Freeze A
FOR epoch = 1 to E_1 DO
Compute residuals: $\hat{r} \leftarrow Y - A X$
Train neural network to minimize $\|\hat{r} - \hat{f}(x; \theta)\|^2$
Update θ using optimizer (SGD/Adam) with learning rate η_θ
END FOR
3. Stage 2 - Joint Fine-tuning:
FOR epoch = 1 to E_2 DO
Compute predictions: $\hat{y} \leftarrow A X + \hat{f}(X; \theta)$
Compute loss: $L \leftarrow \|Y - \hat{y}\|^2 + \lambda(\|A\|_F^2 + \Omega(\theta))$
Compute gradients: $\partial L / \partial A, \partial L / \partial \theta$
Update A: $A \leftarrow A - \eta_A (\partial L / \partial A)$
Update θ : $\theta \leftarrow \theta - \eta_\theta (\partial L / \partial \theta)$ [Note: $\eta_A \ll \eta_\theta$]
Monitor validation loss, apply early stopping if needed
END FOR

Output: Learned A, θ

Hyperparameter Guidance :

- Network widths: 32–256 hidden units (task dependent)
- Hidden layers: 1–3 (typically 2)
- Activation: ReLU or GELU
- Weight decay: 10^{-5} to 10^{-3}
- Dropout rate: 0.1–0.3 (if regularization needed)
- Learning rate for A: $\eta_A \approx 0.001 \times \eta_\theta$ (slow baseline update)

5.3 Model Evaluation Metrics

For Performance Prediction :

- Mean Squared Error (MSE): $MSE = \frac{1}{N} \sum (y_i - \hat{y}_i)^2$
- Mean Absolute Error (MAE): $MAE = \frac{1}{N} \sum |y_i - \hat{y}_i|$
- Inference Time: Measured in milliseconds for practical deployment
- Parameter Efficiency: Total trainable parameters

For Invoice Processing :

- **Precision and Recall:** Per-entity metrics (invoice number, vendor name, total amount, etc.)
- **Character Error Rate (CER):** For OCR components
- **End-to-End Accuracy:** Percentage of invoices with all fields correctly extracted

- **Processing Throughput:** Invoices processed per second
- **Error Rates by Format:** Separate metrics for PDFs, handwritten, scanned documents

5.4 Deployment Architecture

5.4.1 Performance Prediction Deployment

1. **Batch Prediction Server:** For offline analysis (resource planning, capacity modeling)
 - Accept batch workload traces as input
 - Compute baseline predictions: $y_0 = Ax$
 - Inference neural residual: $r = \hat{f}(x; \theta)$
 - Return combined predictions: $y = y_0 + r$
2. **Real-time Prediction API:** For online decision-making (request scheduling, resource allocation)
 - Low-latency inference (target: <1 ms per prediction)
 - Catching of frequently accessed predictions
 - Fallback to baseline if neural network is unavailable
3. **Monitoring and Feedback:** Continuous evaluation against actual system performance
 - Track prediction accuracy over time
 - Detect data drift (workload distribution changes)
 - Trigger periodic retraining when accuracy degrades

5.4.2 Invoice Processing Deployment

1. **Document Ingestion Pipeline:** Accept invoices from multiple sources
 - Email attachments
 - Web upload portal
 - EDI feeds
 - Scanned paper documents via scanning service
2. **OCR Processing:** Parallel processing of documents
 - Use cloud OCR service (e.g., AWS Textract, Google Document AI) or open-source OCR (Tesseract)
 - Queue-based architecture for scalability
 - Retry mechanism for failed conversions
3. **LLM Inference:** Entity extraction via fine-tuned LLM
 - Batch processing for cost efficiency
 - Real-time processing for urgent invoices
 - Catching of repeated vendor patterns

4. **Human Validation Workflow:** Quality assurance and feedback

- Queue suspected errors for human review (e.g., bottom 10% by LLM confidence)
- Track correction patterns to identify systematic issues
- Feed corrections back to improve LLM via fine-tuning or prompt engineering

5. **Database and API:** Downstream integration

- Store extracted entities in relational or vector database
- Expose API for accounting, audit, and analytics systems
- Implement audit trail and compliance logging

5.5 Online Learning and Continuous Improvement

Both systems benefit from online learning and feedback loops :

Performance Models:

- Periodically retrain recent workload data to adapt to changing characteristics
- Optionally adapt A slowly to new resource constraints
- Use MC dropout or ensemble methods for uncertainty estimation

Invoice Processing:

- Feedback from human validation directly improves the LLM via fine-tuning or in-context learning
- Accumulate a growing dataset of validated vendor invoice formats to improve generalization
- Monitor extraction accuracy by vendor and document format to detect systematic issues

6 Limitations and Future Work

6.1 Current Limitations

Digital System Performance Modeling :

- Assumes stationarity in workload patterns; nonstationary workloads (seasonal trends, sudden spikes) require adaptive models
- Neural residuals may overfit with limited data; more data and regularization are needed for robust generalization
- Linear baseline may not capture long-range dependencies or complex interactions
- Cross-system transfer learning remains an open challenge

Invoice Processing :

- OCR performance degrades for low-quality scans or unusual fonts

- LLM-based extraction may struggle with adversarial invoice layouts or intentionally obscured information
- Requires labeled data for fine-tuning; annotation effort is nontrivial for diverse vendor bases
- Performance varies significantly by document format (PDFs > scans > handwritten)

Common Limitations:

- Hybrid approaches require balancing two components; finding optimal split is domain-dependent
- Interpretability of the neural residual remains limited; residual networks remain partially black box
- Both approaches assume adequate training data; few-shot or zero-shot scenarios remain challenging

6.2 Future Directions

Enhanced Performance Modeling :

- Incorporate transformer-based architectures as residual learners, leveraging sequential dependencies in workload traces
- Extend to multi-node distributed architectures where communication bottlenecks dominate
- Add uncertainty-aware prediction using Bayesian neural networks or ensemble methods
- Develop efficient online adaptation strategies for nonstationary workloads

Advanced Invoice Processing :

- Integrate layout-aware model components (e.g., vision transformers) to better understand 2D document structure
- Explore multimodal LLMs that can jointly process image and OCR text, reducing OCR-related errors
- Develop domain-specific LLMs trained on finance and e-commerce invoice corpora
- Implement few-shot adaptation to new vendor types with minimal labeled examples

Unified Research Directions:

- **Explainable Hybrid Models:** Develop interpretability techniques for neural residuals, moving beyond the baseline's transparency
- **Theoretical Analysis:** Provide PAC-learning bounds and generalization guarantees for hybrid decompositions
- **Meta-Learning:** Learn how to decompose problems (baseline selection, residual architecture) across diverse domains
- **Robust Hybrid Models:** Develop hybrid approaches resistant to distribution shift and adversarial inputs

- **Energy Efficiency:** Quantize and compress hybrid models for edge deployment with minimal accuracy loss

7 Conclusion

This report has synthesized two complementary case studies in hybrid artificial intelligence architectures: (1) linear algebra-based digital system performance modeling with neural residuals, and (2) OCR-integrated large language models for automated invoice processing in e-commerce. Both exemplify a powerful design principle: **decompose complex problems into interpretable baselines plus learned residuals.**

The hybrid approach delivers substantial practical benefits :

- **Accuracy:** Near state-of-the-art performance (MSE 0.014 vs. 0.012 for pure NN, 75% error reduction in invoice extraction)
- **Efficiency:** Dramatic reduction in computational cost and parameters (8.5k vs. 200k parameters ; 90% labor reduction)
- **Interpretability:** Baseline remains transparent, enabling diagnosis and debugging
- **Scalability:** Cost-efficient scaling to large workloads and datasets

The report has provided mathematical foundations (error bounds, complexity analysis), implementation guidance (algorithms, hyperparameters, deployment architectures), and empirical validation across two distinct domains. These elements collectively demonstrate the generality and robustness of the hybrid decomposition principle.

Key Takeaways for Practitioners:

1. When faced with complex prediction or extraction tasks, consider decomposing the problem into a cheap, interpretable baseline plus a learned residual.
2. The baseline should encode domain knowledge; the residual should capture phenomena the baseline misses.
3. Hybrid models enable a favorable trade-off between accuracy, interpretability, and computational efficiency.
4. Both components should be jointly optimized, but with careful hyperparameter tuning (e.g., learning rates).
5. Continuous feedback and monitoring are essential for maintaining performance as data distributions change.

As artificial intelligence systems scale to real-world applications, the hybrid paradigm—balancing interpretability, efficiency, and accuracy—represents a pragmatic and effective approach to building systems that are not only accurate but also trustworthy, efficient, and maintainable.

8 Conflict of Interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

9 Author Contributions

Goutam Parashuram Gotur handled the case study, including its design, execution, and analysis, while Dr. E. Saravana Kumar guided the case study by providing supervision, methodological direction, and critical insights. Both authors reviewed the manuscript thoroughly and approved the final version.

10 Funding

No external funding was received for this research.

11 Acknowledgments

We gratefully acknowledge the contributions of the Department of Computer Science and Engineering at The Oxford College of Engineering. We thank colleagues and reviewers for valuable feedback that improved this manuscript.

12 Data Availability Statement

Synthetic performance prediction datasets and invoice processing pipelines are available upon request from the corresponding authors. Real-world data cannot be released due to proprietary and confidentiality constraints.

13 References

- [1] Patel, D., & Pandit, H. B. (2024). Case study: Centralising diverse e-commerce invoices using invoice LLM model. *Scientific Research Journal of Science, Engineering and Technology*, 2(2), 79–82. Retrieved from <https://isrdo.org/journal/SRJSET/currentissue/case-study-centralizing-diverse-e-commerce-invoices-using-invoice-llm-model>
- [2] Sankaran, A., Alashtiy, N. A., & Psarras, C. (2022). Benchmarking the linear algebra awareness of TensorFlow and PyTorch. RWTH Aachen University. Retrieved from <https://arxiv.org/pdf/2202.09888>
- [3] Pudukkottai, et al. (2021). Linear algebraic methods in neural networks. *International Journal of Engineering Research & Technology*, 12(1), 035.

- [4] Baggag, A., & Saad, Y. (2023). Deep learning, transformers and graph neural networks: A linear algebra perspective. Qatar Computing Research Institute & University of Minnesota. Retrieved from <https://www-users.cse.umn.edu/~saad/PDF/nla4dnns.pdf>
- [5] Desai, D., Jain, A., Naik, D., Panchal, N., & Sawant, D. (2021). Invoice processing using RPA & AI. *SSRN Electronic Journal*. <https://doi.org/10.2139/ssrn.3852575>
- [6] Baviskar, D., Ahirrao, S., Potdar, V., & Kotecha, K. (2021). Efficient automated processing of the unstructured documents using artificial intelligence: A systematic literature review and future directions. *IEEE Access*, 9, 72894–72936. <https://doi.org/10.1109/ACCESS.2021.3072900>
- [7] Saout, T., Lardeux, F., & Saubion, F. (2024). An overview of data extraction from invoices. *IEEE Access*. <https://doi.org/10.1109/ACCESS.2024.3360528>
- [8] Bardvall, M., & Hassle, I. (2024). Automating invoice recognition: A comparative study of large language models and OCR/ML technologies.
- [9] Daqqah, B. H. (2024). Leveraging large language models (LLMs) for automated extraction and processing of complex ordering forms. *Doctoral Dissertation*, Massachusetts Institute of Technology.