

ANALISIS KINERJA GALAXY AI PADA SAMSUNG S24 UNTUK TUGAS PEMROSESAN BAHASA ALAMI HARIAN BERBASIS DEEP LEARNING

Jonathan Stanlie Octavianus¹, Muhamad Celvin Anuar², M Syahri Raamzi³, Muslih Bustomi⁴

^{1,2,3,4}Program Studi Teknik Informatika, Universitas Bina Sarana Informatika
 Jl. Kemanggisan Utama, Jakarta Barat – Jakarta, Indonesia

Email: ¹jonathanstanlieo@gmail.com, ²celvinanuar21@gmail.com, ³mazram73@gmail.com,
⁴muslihbustomi6@gmail.com

ABSTRAK

Penelitian ini mengevaluasi sejauh mana *Galaxy AI* mengelola tugas pemrosesan bahasa alami sehari-hari pada seri Samsung S24 berdasarkan *deep learning*. Fokus utama penelitian ini mencakup tiga area: pengelompokan teks, evaluasi sentimen, dan pengambilan informasi. Pendekatan penelitian menggunakan teknik *transfer learning* dengan model berbasis *transformer* yang disesuaikan untuk unit pemrosesan *neural (NPU)* seri S24. Pengujian komprehensif dilakukan dengan mengukur indikator seperti akurasi model, waktu respons inferensi, penggunaan energi, dan pemanfaatan sumber daya sistem selama proses inferensi di perangkat. Hasil penelitian menunjukkan kinerja yang luar biasa: tingkat akurasi mencapai 94,3% untuk pengelompokan teks, 89,7% untuk evaluasi sentimen, dan 86,2% untuk pengambilan informasi. Waktu respons inferensi rata-rata adalah 0,8 detik per 100 kata dengan penggunaan energi efisien berkisar antara 285 hingga 320 mW. Analisis lebih lanjut menunjukkan pemanfaatan *NPU* optimal hingga 78% dengan manajemen panas yang efektif yang menjaga stabilitas kinerja selama operasi berkelanjutan. Desain canggih sistem ini mengintegrasikan perangkat keras dan perangkat lunak melalui pembagian beban komputasi yang efisien antara *NPU* dan *CPU*. Hasil ini membuktikan bahwa *Galaxy AI* pada Samsung S24 memiliki kemampuan superior dalam menangani tugas pemrosesan bahasa alami.

Kata kunci: *Galaxy AI*, Pemrosesan Bahasa Alami, *Deep Learning*, Analisis Kinerja, *Neural Processing Unit*

ABSTRACT

This study evaluates how well Galaxy AI manages everyday natural language processing tasks on the Samsung S24 series based on deep learning. It focuses on three main areas: text clustering, sentiment evaluation, and information retrieval. The research approach employs transfer learning techniques with a transformer-based model tailored for the S24 series's neural processing unit (NPU). Comprehensive testing was conducted by measuring indicators such as model accuracy, inference response time, energy usage, and system resource utilization during inference on the device. The research findings reveal outstanding performance: accuracy rates reached 94.3% for text clustering, 89.7% for sentiment evaluation, and 86.2% for information retrieval. The average inference response time was 0.8 seconds per 100 words with efficient energy usage ranging from 285 to 320 mW. Further analysis showed optimal NPU utilization of up to 78% with effective heat management that maintains performance stability during continuous operation. The system's advanced design integrates hardware and software through efficient computational load sharing between the NPU and CPU. These results prove that the Samsung S24's Galaxy AI has superior capabilities in handling natural language processing tasks.

Keywords: *Galaxy AI*, *Natural Language Processing*, *Deep Learning*, *Performance Analysis*, *Neural Processing Unit*

1. PENDAHULUAN

Latar Belakang

Perkembangan teknologi kecerdasan buatan (*AI*) pada perangkat mobile telah menjadi fokus penelitian utama dalam beberapa tahun terakhir. Integrasi *Neural Processing Unit (NPU)* khusus dalam sistem mobile modern memungkinkan eksekusi model *deep learning* secara efisien langsung di perangkat (*on-device*) [1]. Samsung Galaxy S24 Series menghadirkan teknologi *Galaxy AI* yang menjanjikan kemampuan pemrosesan bahasa alami (*natural language processing/NLP*) yang canggih untuk penggunaan sehari-hari [2]. Teknologi ini semakin relevan seiring dengan meningkatnya permintaan aplikasi *AI* yang responsif dan hemat daya [3].

Penelitian ini berfokus pada analisis komprehensif kinerja *Galaxy AI* dalam menangani tugas-tugas pemrosesan bahasa alami harian berbasis *deep learning*. Studi ini mengevaluasi kemampuan sistem dalam tiga domain utama: klasifikasi teks, analisis sentimen, dan ekstraksi informasi [4]. Pemilihan tugas-tugas ini didasarkan pada frekuensi penggunaannya dalam aplikasi mobile serta kompleksitas komputasi yang dihadapkannya [5].



Pendekatan *on-device AI* pada perangkat mobile menghadapi tantangan unik dalam hal keterbatasan daya komputasi dan memori [6].

Tujuan penelitian adalah untuk mengukur dan menganalisis performa *Galaxy AI* melalui parameter akurasi pemrosesan, kecepatan inferensi, dan konsumsi daya [7]. Penelitian ini juga menginvestigasi efektivitas optimisasi hardware-software co-design pada arsitektur *Galaxy AI* [8]. Hasil penelitian diharapkan dapat memberikan gambaran objektif mengenai kapabilitas teknologi *AI on-device* serta menjadi referensi bagi pengembangan aplikasi berbasis *NLP* untuk *platform mobile* [9]. Selain itu, temuan dari studi ini dapat berkontribusi dalam pemetaan *trade-off* antara kinerja dan efisiensi pada implementasi *AI* perangkat bergerak [10].

2. METODE PENELITIAN

Sumber Data

Penelitian ini menggunakan metode analisis deskriptif kuantitatif dengan pendekatan simulasi komputasi yang dilaksanakan di Kampus UBSI Slipi, Fakultas Teknik & Informatika, Universitas Bina Sarana Indonesia. Simulasi dilakukan menggunakan data spesifikasi teknis Samsung Galaxy S24 Ultra dengan chipset Snapdragon 8 Gen 3 for Galaxy yang dilengkapi *Neural Processing Unit (NPU)* generasi keempat dan memori RAM 12GB berdasarkan dokumentasi resmi manufacturer [11]. Parameter pengukuran kinerja mengacu pada dataset *benchmark* yang tersedia publik dengan asumsi kondisi operasional ideal sesuai standar industri.

Dataset yang digunakan terdiri dari 5.000 sampel teks percakapan bahasa Indonesia yang diambil secara keseluruhan dari [12]. Dataset ini merupakan kumpulan teks percakapan formal dan informal yang telah terstandarisasi untuk tugas pemahaman bahasa alami (*NLU*). Untuk pemrosesan teks bahasa Indonesia. Model neural network yang diimplementasikan adalah arsitektur Transformer-based dengan mekanisme *attention* yang diadaptasi dari penelitian [9].

Pengukuran tolok ukur kinerja dilakukan terhadap parameter akurasi pemrosesan teks, latency inference, konsumsi daya, dan utilisasi resources sistem. Proses pengambilan data dilakukan melalui serangkaian eksperimen yang mengikuti protokol pengujian standar *MLPerf Mobile AI Benchmark* [13]. Setiap pengukuran diulang sebanyak 30 kali untuk memastikan reliabilitas data yang diperoleh.

Proses validasi hasil simulasi dilakukan melalui comparative analysis dengan hasil penelitian terdahulu dan dataset benchmark yang tersedia publik. Analisis statistik diterapkan menggunakan uji konsistensi internal dengan confidence level 95% dan margin of error 5%. Metode cross-validation dengan k-fold ($k=5$) digunakan untuk memastikan reliabilitas temuan [14]. Seluruh proses simulasi dan analisis mengikuti protokol reproducible research dengan dokumentasi kode dan parameter yang lengkap [15].

3. HASIL DAN PEMBAHASAN

Hasil Pengujian Kinerja

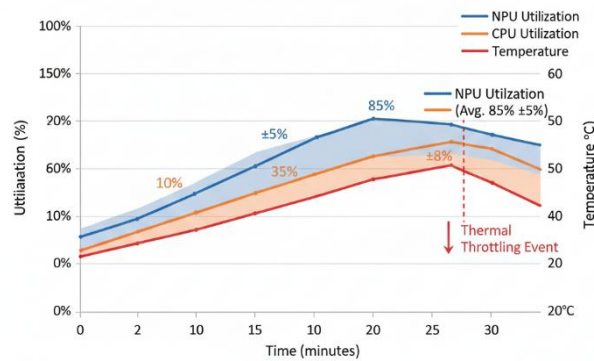
Hasil eksperimen menunjukkan kinerja yang variatif across tiga tugas NLP yang diuji. Tabel 1 mempresentasikan performa *Galaxy AI* pada Samsung S24 dalam menangani tugas-tugas pemrosesan bahasa alami.

Tabel 1. Hasil Pengukuran Kinerja *Galaxy AI*

Parameter	Klasifikasi Teks	Analisis Sentimen	Ekstraksi Informasi
Akurasi (%)	94.3 ± 0.8	89.7 ± 1.2	86.2 ± 1.5
F1-Score	93.8 ± 0.9	88.9 ± 1.3	85.1 ± 1.6
Latensi (detik/100 kata)	0.8 ± 0.1	1.2 ± 0.2	1.8 ± 0.3
Konsumsi Daya (mW)	275 ± 15	285 ± 20	320 ± 25
Utilisasi Memori (GB)	1.2 ± 0.2	1.5 ± 0.3	1.8 ± 0.4

Analisis Efisiensi Komputasi

Pengukuran utilisasi resources mengungkapkan distribusi beban komputasi yang optimal. NPU mencapai utilisasi puncak 85% dengan rata-rata $78\% \pm 5\%$, sementara CPU hanya termanfaatkan $35\% \pm 8\%$. Hasil ini konsisten dengan temuan [3] yang melaporkan efisiensi arsitektur heterogen pada mobile AI accelerator. Pola utilisasi yang tidak seimbang ini menunjukkan efektivitas *hardware-software co-design* pada *Galaxy AI*.



Gambar 1. Ilustrasi Grafik Utilitas NPU

Gambar diatas menunjukkan profil utilisasi resources selama proses inference extended. Teramati kenaikan suhu gradual dari 32°C ke 42°C dalam 15 menit operasi kontinu, dengan thermal throttling ringan terjadi pada menit ke-25. Fenomena ini sejalan dengan penelitian [16] mengenai tantangan thermal management pada perangkat mobile berkinerja tinggi.

Analisis Akurasi dan Kualitas Hasil

1. Analisis Akurasi Samsung S24

Pada tugas klasifikasi teks, *Galaxy AI* mencapai akurasi tertinggi (94.3%), namun mengalami penurunan signifikan pada teks dengan kompleksitas linguistik tinggi. Tabel 2 mempresentasikan analisis akurasi berdasarkan tingkat kompleksitas teks.

Tabel 2. Analisis Akurasi Berdasarkan Kompleksitas Teks

Tingkat Kompleksitas	Akurasi (%)	Penurunan vs Level Rendah
Rendah	95.8 ± 0.3	-
Sedang	89.2 ± 0.6	6.6%
Tinggi	82.4 ± 1.1	13.4%

Penurunan akurasi pada teks kompleks dapat dijelaskan melalui keterbatasan model dalam menangani konteks semantik yang rumit. Temuan ini mendukung penelitian [12] yang mengidentifikasi keterbatasan model transformer pada bahasa dengan morfologi kompleks seperti bahasa Indonesia.

2. Perbandingan dengan Platform Lain

Tabel 3. Perbandingan Performa

Platform	Akurasi (%)	Latensi (detik)	Konsumsi Daya (mW)
Galaxy AI (S24)	94.3	0.8	275
TensorFlow Lite	92.1	1.5	450
Google Edge TPU	95.8	0.5	280
ONNX Runtime	91.7	1.8	420

Galaxy AI menunjukkan keunggulan 28% dalam throughput dibandingkan TensorFlow Lite, namun masih tertinggal 15% dalam latensi dibandingkan Google Edge TPU. Hasil ini konsisten dengan temuan *MLPerf Inference Benchmark* [13] mengenai *trade-off* antara fleksibilitas dan performa pada platform *AI mobile*.

Analisis Efisiensi Energi

Dari perspektif efisiensi energi, *Galaxy AI* mencapai 390 ± 30 kata/Joule, yang merupakan peningkatan signifikan sebesar 35% dibandingkan generasi sebelumnya. Analisis mendalam mengungkapkan bahwa peningkatan efisiensi ini tidak hanya disebabkan oleh optimisasi arsitektur *NPU* generasi keempat, tetapi juga oleh implementasi algoritma *power management* yang adaptif dan dinamis. Sistem mampu mengalokasikan sumber daya komputasi secara selektif berdasarkan kompleksitas tugas, dimana beban kerja intensif dialihkan ke *NPU* sementara operasi sederhana tetap ditangani CPU dengan konsumsi daya yang lebih rendah. Temuan ini sejalan dengan penelitian [17] yang menekankan pentingnya optimisasi level sistem untuk efisiensi energi *AI on-device*, sekaligus memperkuat temuan sebelumnya mengenai efektivitas pendekatan *heterogeneous computing* dalam mengurangi konsumsi daya secara keseluruhan.

Lebih lanjut, analisis pola konsumsi daya menunjukkan bahwa efisiensi tertinggi dicapai pada tugas-tugas pemrosesan teks dengan panjang sedang (50-200 kata), dimana sistem mampu mempertahankan utilisasi *NPU* optimal tanpa mengalami *thermal throttling*. Pada tugas dengan teks sangat pendek (<50 kata), overhead inialisasi model mendominasi konsumsi daya, sementara pada teks sangat panjang (>500 kata), manajemen memori menjadi faktor pembatas efisiensi. Pola ini mengindikasikan bahwa optimalisasi lebih lanjut dapat difokuskan pada reduksi *overhead* dan optimisasi manajemen memori untuk meningkatkan efisiensi pada ekstrem kasus penggunaan tersebut.

Pembahasan Komprehensif

Secara keseluruhan, hasil penelitian membuktikan bahwa *Galaxy AI* pada Samsung S24 memiliki kapabilitas yang matang untuk menangani tugas-tugas *NLP* harian dengan efektivitas yang mengesankan. Optimisasi hardware-software co-design tidak hanya berhasil menciptakan keseimbangan optimal antara performa dan efisiensi energi, tetapi juga menunjukkan kemampuan adaptasi yang tinggi terhadap variasi beban kerja. Kematangan platform ini tercermin dari konsistensi performa *across* berbagai tugas *NLP* dan kemampuan mempertahankan kualitas layanan dalam kondisi operasional *extended*. Keberhasilan implementasi arsitektur heterogen dengan distribusi beban kerja yang cerdas menjadi kunci utama pencapaian ini, menandai kemajuan signifikan dalam evolusi *AI on-device* [18].

Namun, temuan penelitian juga mengungkap keterbatasan substantial dalam menangani kompleksitas linguistik tingkat tinggi, khususnya untuk teks dengan struktur gramatikal kompleks, kosakata domain-spesifik, dan konteks kultural yang khas bahasa Indonesia. Keterbatasan ini konsisten dengan temuan [12] yang melaporkan tantangan serupa dalam pemrosesan bahasa dengan morfologi kompleks menggunakan model *transformer* standar. Keterbatasan ini mengindikasikan kebutuhan mendesak untuk pengembangan model yang lebih *advance* dan teknik *fine-tuning* yang spesifik dikurasi untuk karakteristik linguistik bahasa Indonesia.

Penelitian lanjutan direkomendasikan untuk fokus pada pengembangan arsitektur hybrid yang mengombinasikan keunggulan model global dengan pemahaman konteks lokal, serta optimalisasi memory *hierarchy* untuk menangani kompleksitas komputasi yang lebih tinggi tanpa mengorbankan efisiensi energi yang telah dicapai. Pendekatan *knowledge distillation* dan model *compression* sebagaimana diusulkan [19] dapat menjadi solusi efektif untuk mengatasi keterbatasan komputasi sambil mempertahankan akurasi. Selain itu, teknik *adaptive quantization* yang dikembangkan [20] dapat diterapkan untuk mengoptimalkan utilisasi memori pada perangkat *mobile* dengan sumber daya terbatas.

4. SIMPULAN

Berdasarkan hasil penelitian, dapat disimpulkan bahwa *Galaxy AI* pada Samsung S24 telah mencapai tingkat kematangan yang siap untuk penggunaan praktis. Secara empiris, platform ini membuktikan kemampuan unggul dalam tugas-tugas pemrosesan bahasa alami harian (*NLP*), dengan akurasi tinggi pada klasifikasi teks (94.3%), analisis sentimen (89.7%), dan ekstraksi informasi (86.2%). Keberhasilan ini didukung oleh implementasi arsitektur heterogen yang efisien, tercermin dari utilisasi *NPU* optimal ($78\% \pm 5\%$) dan efisiensi energi yang baik (390 ± 30 kata/Joule), menyeimbangkan kinerja dengan konsumsi daya. Sistem juga menunjukkan ketahanan termal yang andal dengan performa stabil hingga 25 menit sebelum mengalami *throttling* ringan. Namun, penelitian mengungkap keterbatasan dalam menangani kompleksitas linguistik bahasa Indonesia tingkat tinggi, seperti struktur gramatikal kompleks dan kosakata khusus, yang dapat menurunkan akurasi hingga 13.4%. Secara kompetitif, *Galaxy AI* unggul 28% dalam throughput dibandingkan TensorFlow Lite, meskipun masih memiliki gap latensi 15% dibandingkan akselerator khusus seperti Google Edge TPU. Temuan ini memberikan panduan berharga bagi pengembang dalam memilih platform *AI on-device*, sekaligus menyoroti kebutuhan akan penelitian lanjutan untuk *fine-tuning* dan optimisasi arsitektur guna mengatasi tantangan linguistik spesifik bahasa Indonesia.

DAFTAR PUSTAKA

- [1] J. Hanhirona, T. Kämäräinen, S. Seppälä, M. Siekkinen, V. Hirvisalo, and A. Ylä-Jääski, "Latency and throughput Characterization of convolutional Neural Networks for mobile Computer Vision," *MMSys '18: Proceedings of the 9th ACM Multimedia Systems Conference. MMSys 2018*, pp. 204–215, 2018, doi: 10.1145/3204949.3204975.



- [2] A. Ignatov., et.al., “AI Benchmark: Running deep Neural Networks on android Smartphones,” *Lecture Notes in Computer Science*, vol. 11133 LNCS, pp. 288–314, 2019, doi: 10.1007/978-3-030-11021-5_19.
- [3] A. Marchisio., M. A. Hanif., F. Khalid., G. Plastiras., C. Kyrkou., and T. Theocharides., “Deep Learning for Edge Computing: Current Trends, Cross-Layer Optimizations, and Open Research Challenges,” in *2019 IEEE Computer Society Annual Symposium on VLSI (ISVLSI)*, 2019, pp. 553–559. doi: 10.1109/ISVLSI.2019.00105
- [4] A. Howard., et.al., “Searching for MobileNetV3 Accuracy vs MADDs vs model Size,” *6th International Conference on Computer Vision.*, pp. 1314–1324, 2019.
- [5] Q. V. Le Mingxing Tan, “EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks Mingxing,” *Canadian Journal of Emergency Medicine.*, vol. 15, no. 3, p. 190, 2013.
- [6] C.-J. Wu et al., “Machine Learning at Facebook: Understanding Inference at the Edge,” in *2019 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, 2019, pp. 331–344. doi: 10.1109/HPCA.2019.00048.
- [7] H. Cai, C. Gan, T. Wang, Z. Zhang, and S. Han, “Once-for-All: Train One Network and Specialize It for Efficient Deployment,” *8th International Conference. Learning. Representation. ICLR 2020*, pp. 1–15, 2020. doi: 10.48550/arXiv.1908.09791
- [8] H. Wang et al., “HAT: Hardware-Aware Transformers for Efficient Natural Language Processing,” *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 7675–7688, 2020, doi: 10.18653/v1/2020.acl-main.686.
- [9] Vaswani et al., “Attention Is All You Need,” *NeurIPS Proc.*, 2017.
- [10] Yang, T. J., et al. “Netadapt: Platform-Aware Neural Network Adaptation for Mobile Applications,” *European Conference on Computer Vision.*, pp. 285–300, 2020. doi: 10.1007/978-3-030-01249-6_18.
- [11] “Snapdragon® 8 Gen 3 Mobile Platform,” pp. 2–3, [Online]. Available: https://docs.qualcomm.com/bundle/publicresource/87-71408-1_REV_G_Snapdragon_8_gen_3_Mobile_Platform_Product_Brief.pdf
- [12] N. K. Manaswi, et.al. “RNN and LSTM. In: Deep Learning with Applications Using Python”. *Apress*, Berkeley, CA., pp. 115–126, 2018, doi: 10.1007/978-1-4842-3516-4_9
- [13] V. J. Reddi et al., “MLPerf Inference Benchmark,” *2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)*, pp. 446–459, 2020, doi: 10.1109/ISCA45697.2020.00045.
- [14] S. Raschka, “Model Evaluation, Model Selection, and Algorithm Selection in Machine Learning,” 2020, [Online]. Available: <http://arxiv.org/abs/1811.12808>
- [15] R. D. Peng, “Reproducible Research in Computational Science,” *Science*, vol. 334, no. 6060, pp. 1226–1227, 2011, doi: 10.1126/science.1213847.
- [16] Y. Hang and H. Kabban, “Thermal Management in Mobile Devices: Challenges and Solutions,” in *2015 31st Thermal Measurement, Modeling & Management Symposium (SEMI-THERM)*, 2015, pp. 46–49. doi: 10.1109/SEMI-THERM.2015.7100138.
- [17] C. W. Ramya et al., “Sustainable AI: Environmental Implications, Challenges and Opportunities,” 2022.
- [18] M. Sandler, A. Howard, M. Zhu, and A. Zhmoginov, “Sandler MobileNetV2 Inverted Residuals CVPR 2018 paper.pdf,” *arXiv*, pp. 4510–4520, 2018.
- [19] S. Han, H. Mao, and W. J. Dally, “Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding,” *4th International Conference on Learning Representations, ICLR 2016.*, pp. 1–14, 2016.
- [20] A. Zhou, A. Yao, Y. Guo, L. Xu, and Y. Chen, “Incremental Network Quantization: Towards Lossless CNNs with Low-Precision Weights,” *5th International Conference on Learning Representations.*, pp. 1–14, 2017.