



**IMPLEMENTASI MODEL *VISUAL QUESTION ANSWERING*
MENGUNAKAN *VISION TRANSFORMER* DAN
EFFICIENTNET-V2 DENGAN BERT**

AHMED NIZHAN HAIKAL

2110511022

**INFORMATIKA
FAKULTAS ILMU KOMPUTER
UNIVERSITAS PEMBANGUNAN NASIONAL "VETERAN"
JAKARTA
JAKARTA
2025**



**IMPLEMENTASI MODEL *VISUAL QUESTION ANSWERING*
MENGUNAKAN *VISION TRANSFORMER* DAN
EFFICIENTNET-V2 DENGAN BERT**

SKRIPSI

**Diajukan Sebagai Salah Satu Syarat Untuk Memperoleh Gelar Sarjana
Komputer**

AHMED NIZHAN HAIKAL

2110511022

INFORMATIKA

FAKULTAS ILMU KOMPUTER

UNIVERSITAS PEMBANGUNAN NASIONAL “VETERAN” JAKARTA

JAKARTA

2025

PERNYATAAN ORISINALITAS

Tugas akhir ini adalah hasil karya sendiri dan semua sumber baik yang dikutip maupun yang dirujuk telah saya nyatakan dengan benar.

Nama : Ahmed Nizhan Haikal

NIM : 2110511022

Tanggal : 29 Juni 2025

Bilamana di kemudian hari ditemukan ketidaksesuaian dengan pernyataan saya ini, maka saya bersedia dituntut dan diproses sesuai dengan ketentuan yang berlaku.

Bekasi, 29 Juni 2025

Yang Menyatakan



Ahmed Nizhan Haikal

**PERNYATAAN PERSETUJUAN PUBLIKASI TUGAS AKHIR UNTUK
KEPENTINGAN AKADEMIS**

Sebagai civitas akademika Universitas Pembangunan Nasional Veteran Jakarta,
saya yang bertanda tangan di bawah ini:

Nama : Ahmed Nizhan Haikal
NIM : 2110511022
Fakultas : Ilmu Komputer
Program Studi : S-1 Informatika

Demi pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada Universitas Pembangunan Nasional Veteran Jakarta Hak Bebas Royalti Non eksklusif (Non - exclusive Royalty Free Right) atas skripsi saya yang berjudul:

**IMPLEMENTASI MODEL *VISUAL QUESTION ANSWERING*
MENGUNAKAN *VISION TRANSFORMER* DAN *EFFICIENTNET-V2*
DENGAN BERT**

Dengan Hak Bebas Royalti ini Universitas Pembangunan Nasional Veteran Jakarta berhak menyimpan, mengalih media/memformatkan, mengelola dalam bentuk pangkalan data (basis data), merawat dan mempublikasikan skripsi saya selama tetap mencantumkan nama saya sebagai penulis dan sebagai pemilik Hak Cipta. Demikian pernyataan ini saya buat dengan sebenarnya.

Dibuat di: Bekasi
Pada tanggal: 29 Juni 2025
Yang Menyatakan



Ahmed Nizhan Haikal

LEMBAR PENGESAHAN

Judul : Implementasi Model *Visual Question Answering* Menggunakan *Vision Transformer* dan *EfficientNet-V2* dengan BERT
Nama : Ahmed Nizhan Haikal
NIM : 2110511022
Program Studi : S1 Informatika

Disetujui oleh :

Penguji 1:

Neny Rosmawarni, S.Kom, M.Kom.

Penguji 2:

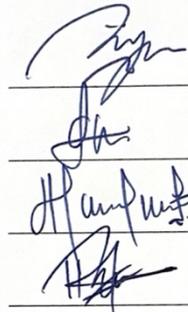
Nurul Afifah Arifuddin, S. Pd., M.T.

Pembimbing 1:

Ridwan Raafi'udin, S.Kom., M.Kom.

Pembimbing 2:

Muhammad Adrezo, S.Kom., M.Sc.



Diketahui oleh:

Koordinator Program Studi:

Dr. Widya Cholil, M.I.T

NIP. 221112080

Dekan Fakultas Ilmu Komputer:

Prof. Dr. Ir. Supriyanto, M.Sc., IPM.

NIP. 197605082003121002



Tanggal Ujian Tugas Akhir:

02 Juni 2025

IMPLEMENTASI MODEL *VISUAL QUESTION ANSWERING* MENGUNAKAN *VISION TRANSFORMER* DAN *EFFICIENTNET-V2* DENGAN BERT

Ahmed Nizhan Haikal

ABSTRAK

Visual Question Answering (VQA) merupakan tugas *multimodal* yang memadukan visi komputer dan pemrosesan bahasa alami untuk menjawab pertanyaan berdasarkan citra. Penelitian ini mengeksplorasi pengaruh arsitektur visual (*EfficientNet V2* dan *Vision Transformer/ViT*) terhadap performa VQA, dengan BERT digunakan sebagai *backbone* teks untuk pemahaman konteks semantik *bidirectional*. Penelitian ini menggunakan *dataset* DT-VQA dan menerapkan strategi *thresholding* untuk mengatasi ketidakseimbangan label. Hasil menunjukkan bahwa model ViT-BERT dengan *fine-tuning* mencapai performa terbaik di antara varian ViT-BERT, mencatat akurasi 53,19% dan ANLS 61,95% pada *threshold* 75. Performa terbaik dicatat oleh model *EfficientNet V2*-BERT dengan *transfer learning*, dengan akurasi sebesar 56,3% dan ANLS 62,9%. Temuan ini menggarisbawahi pentingnya penyesuaian arsitektur dan strategi pelatihan untuk optimasi performa pada karakteristik data.

Kata kunci: BERT, *EfficientNet V2*, Pemrosesan *Multimodal* , *Visual Question Answering*, *Vision Transformer*

IMPLEMENTASI MODEL *VISUAL QUESTION ANSWERING* MENGUNAKAN *VISION TRANSFORMER* DAN *EFFICIENTNET-V2* DENGAN BERT

Ahmed Nizhan Haikal

ABSTRACT

Visual Question Answering (VQA) is a multimodal task that integrates computer vision and natural language processing to answer image-based questions. This research explores the impact of visual architectures (EfficientNet V2 and Vision Transformer/ViT) on VQA performance, utilizing BERT as the text backbone for bidirectional semantic context understanding. We employed the DT-VQA dataset and implemented a thresholding strategy to address label imbalance. Results indicate that the ViT-BERT model with fine-tuning achieved the best performance among ViT-BERT variants, recording an accuracy of 53.19% and an ANLS of 61.95% at a threshold of 75. Overall, the highest performance was achieved by the EfficientNet V2-BERT model with transfer learning, with an accuracy of 56.3% and an ANLS of 62.9%. These findings underscore the importance of tailoring architectural choices and training strategies for optimal performance given specific data characteristics.

Keywords: BERT, EfficientNet V2, Multimodal Processing, Visual Question Answering, Vision Transformer

Kata Pengantar

Puji syukur penulis panjatkan kepada Tuhan Yang Maha Esa karena atas rahmat dan karunia-Nya, penulis dapat menyelesaikan skripsi yang berjudul "Implementasi Model *Visual Question Answering* Menggunakan *Vision Transformer* dan *EfficientNet-V2* Dengan BERT". Skripsi ini disusun sebagai salah satu syarat untuk menyelesaikan program studi sarjana pada Program Studi Informatika, Fakultas Ilmu Komputer, Universitas Pembangunan Nasional "Veteran" Jakarta. Dalam penulisan skripsi ini, penulis menyadari bahwa banyak bantuan, dukungan, dan bimbingan yang penulis terima dari berbagai pihak. Oleh karena itu, pada kesempatan ini penulis ingin menyampaikan rasa terima kasih yang sebesar-besarnya kepada:

1. Orang tua dan keluarga yang senantiasa memberikan dukungan, semangat, arahan, dan doa selama pengerjaan skripsi ini.
2. Bapak Prof. Dr. Ir. Supriyanto, ST., M.Sc., IPM selaku Dekan Fakultas Ilmu Komputer Universitas Pembangunan Nasional "Veteran" Jakarta yang telah memberikan izin dan kesempatan untuk menjalankan penelitian ini.
3. Ibu Dr. Widya Cholil, S.Kom., M.I.T selaku Kepala Program Studi S1-Informatika Fakultas Ilmu Komputer Universitas Pembangunan Nasional "Veteran" Jakarta yang telah memberikan arahan dan dukungan selama penyusunan proposal ini.
4. Bapak Ridwan Raafi'udin, M.Kom., selaku dosen pembimbing pertama, yang telah memberikan arahan, saran, dan bimbingan selama proses penyusunan proposal ini.
5. Bapak Muhammad Adrezo, S.Kom., M.Sc. selaku dosen pembimbing kedua yang telah memberikan masukan, saran, serta dukungan selama penulisan proposal ini.
6. Teman-teman yang telah memberikan semangat dan bantuan dalam menyelesaikan skripsi.

Serta semua pihak yang tidak dapat penulis sebutkan satu per satu, yang telah memberikan kontribusi dalam menyelesaikan skripsi.

Daftar Isi

PERNYATAAN ORISINALITAS.....	III
PERNYATAAN PERSETUJUAN PUBLIKASI TUGAS AKHIR UNTUK KEPENTINGAN AKADEMIS.....	IV
LEMBAR PENGESAHAN	I
KATA PENGANTAR.....	IV
DAFTAR ISI	V
DAFTAR GAMBAR	VIII
DAFTAR TABEL.....	IX
DAFTAR RUMUS.....	X
DAFTAR LAMPIRAN	XI
BAB 1. PENDAHULUAN.....	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah	3
1.3 Batasan Masalah.....	4
1.4 Tujuan dan Manfaat Penelitian	5
1.5 Sistematika Penulisan.....	5
BAB 2. TINJAUAN PUSTAKA.....	8
2.1 <i>Machine Learning</i>	8
2.2 <i>Deep Learning</i>	9
2.3 <i>Visual Question Answering</i>	9
2.4 <i>Dataset Dense Text Visual Question Answering</i>	10
2.5 <i>Computer Vision</i>	11
2.6 Teknik Praproses Data.....	12
2.7 <i>Image Recognition</i>	13

2.8	<i>Convolutional Neural Network</i>	14
2.9	<i>EfficientNet V2</i>	16
2.10	<i>Natural Language Processing</i>	17
2.11	<i>Transformer</i>	18
2.12	<i>Vision Transformer</i>	22
2.13	BERT.....	23
2.14	<i>Transfer Learning</i>	25
2.15	<i>Multi-GPU Training</i>	25
2.16	<i>Mixed Precision Training</i>	26
2.17	Metriks Evaluasi Model	27
2.17.1	Akurasi	27
2.17.2	<i>AccuracyANLS (AccANLS)</i>	28
2.18	Penelitian Terdahulu.....	29
BAB 3.	METODE PENELITIAN.....	32
3.1	Tahapan Penelitian	32
3.1.1	Identifikasi Masalah	33
3.1.2	Studi Pustaka.....	34
3.1.3	Pengumpulan Data	34
3.1.4	Praproses Data.....	34
3.1.5	Pembuatan Model.....	37
3.1.6	Pelatihan Model	40
3.1.7	Analisis dan Evaluasi Model.....	41
3.1.8	Implementasi dan Integrasi GUI dengan Model	42
3.2	Alat dan Bahan Penelitian.....	43
3.2.1	Perangkat Keras	43
3.2.2	Perangkat Lunak.....	44

3.3	Jadwal Penelitian.....	45
BAB 4.	HASIL DAN PEMBAHASAN.....	46
4.1	Hasil dan Rekomendasi.....	46
4.1.1	Hasil Pengumpulan Data.....	46
4.1.2	Praproses Data.....	51
4.1.3	Praproses Teks.....	57
4.1.4	Pembuatan Model.....	62
4.1.5	Hasil Pelatihan Model.....	69
4.1.6	Analisis Hasil Pelatihan Model.....	73
4.1.7	Implementasi <i>Graphical User Interface</i>	76
BAB 5.	PENUTUP.....	80
5.1	Kesimpulan	80
5.2	Saran.....	80
	DAFTAR PUSTAKA	82
	LAMPIRAN.....	86

Daftar Gambar

Gambar 2.1 Visualisasi <i>Kernel</i> Konvolusi (Géron 2023)	15
Gambar 2.2 Arsitektur <i>Fused-MBConv</i> (Tan dan Le 2021).....	17
Gambar 2.3 Arsitektur transformer (Guo <i>et al.</i> 2022).....	19
Gambar 2.4 Arsitektur <i>Scaled Dot-Product Attention</i> dan <i>Multi-Head Attention</i> (Guo <i>et al.</i> 2022)	21
Gambar 2.5 Arsitektur vision transformer (Guo <i>et al.</i> 2022).....	22
Gambar 2.6 Arsitektur Model BERT (Saputra <i>et al.</i> 2024).....	24
Gambar 3.1 Tahapan Penelitian	33
Gambar 3.2 Tahapan Praproses Citra	34
Gambar 3.3 Tahapan Praproses Teks	36
Gambar 3.4 Arsitektur Model dengan <i>Backbone</i> ViT dan BERT	37
Gambar 3.5 Arsitektur Model dengan <i>Backbone EfficientNet V2</i> dan BERT	39
Gambar 4.1 Contoh Format Anotasi <i>JSON Dataset</i> DT-VQA.....	47
Gambar 4.2 Contoh Kategori Citra dalam <i>Dataset</i> DT-VQA	48
Gambar 4.3 Contoh Interaksi Tanya Jawab pada <i>Dataset</i> DT-VQA	49
Gambar 4.4 Distribusi Jumlah Jawaban terhadap Jumlah Kemunculannya	50
Gambar 4.5 Alur Tahapan pemuatan data	52
Gambar 4.6 Hasil <i>Resize</i> pada Citra.....	54
Gambar 4.7 Kode Tokenisasi Pertanyaan	58
Gambar 4.8 Kode Tokenisasi Jawaban dan Pembentukan Label.....	60
Gambar 4.9 Kode Model VQA dengan <i>Backbone</i> ViT dan BERT	63
Gambar 4.10 Kode <i>Fine Tuning</i> Model	64
Gambar 4.11 Kode Model VQA dengan <i>Backbone EfficientNet V2</i> dan BERT	66
Gambar 4.12 Tampilan Halaman <i>Input</i>	77
Gambar 4.13 Tampilan Aplikasi Setelah Diberikan <i>Input</i>	78
Gambar 4.14 Halaman Hasil	79

Daftar Tabel

Tabel 2.1 Penelitian Terdahulu	29
Tabel 3.1 Jadwal Penelitian	45
Tabel 4.1 Tabel Kemunculan Jawaban Terhadap <i>Threshold</i>	50
Tabel 4.2 Contoh Data Setelah Transformasi ke <i>DataFrame</i>	53
Tabel 4.3 Contoh Konversi Nilai Piksel dan Normalisasi Rentang Nilai	55
Tabel 4.4 Contoh standarisasi pada nilai piksel citra	56
Tabel 4.5 Contoh Normalisasi Teks Jawaban.....	58
Tabel 4.6 Contoh Hasil Tokenisasi Pertanyaan	59
Tabel 4.7 Contoh Hasil Tokenisasi Jawaban dan Pembentukan Label	61
Tabel 4.8 Contoh Perhitungan Metrik <i>Accuracy</i> , ANLS, dan <i>AccuracyANLS</i>	70
Tabel 4.9 Hasil Pelatihan Model dengan <i>Backbone</i> ViT dan BERT	71
Tabel 4.10 Hasil Pelatihan Model dengan <i>Backbone EfficientNet V2</i> dan BERT.	72

Daftar Rumus

Rumus 2.1	19
Rumus 2.2	19
Rumus 2.3	20
Rumus 2.4	20
Rumus 2.5	20
Rumus 2.6	23
Rumus 2.7	23
Rumus 2.8	23
Rumus 2.9	23
Rumus 2.10	28
Rumus 2.11.....	28
Rumus 2.12	28
Rumus 4.1	56