

**EKSPERIMEN MULTI-MODEL DAN
MULTI-PERSONA UNTUK MENGANALISIS
DAMPAK *PERSONA* TERHADAP PENALARAN,
PERILAKU KELUARAN, DAN *HUMAN BIAS*
PADA LARGE LANGUAGE MODEL**

Proposal Tugas Akhir

Oleh

**Abel Apriliani
18222008**



**PROGRAM STUDI SISTEM DAN TEKNOLOGI INFORMASI
SEKOLAH TEKNIK ELEKTRO DAN INFORMATIKA
INSTITUT TEKNOLOGI BANDUNG
November 2025**

LEMBAR PENGESAHAN

EKSPERIMEN MULTI-MODEL DAN MULTI-PERSONA UNTUK MENGANALISIS DAMPAK *PERSONA* TERHADAP PENALARAN, PERILAKU KELUARAN, DAN *HUMAN BIAS* PADA LARGE LANGUAGE MODEL

Proposal Tugas Akhir

Oleh

Abel Apriliani
18222008

Program Studi Sistem dan Teknologi Informasi
Sekolah Teknik Elektro dan Informatika
Institut Teknologi Bandung

Proposal Tugas Akhir ini telah disetujui dan disahkan
di Bandung, pada tanggal 16 November 2025

Pembimbing 1

Pembimbing 2

Dr. Eng. Ayu Purwarianti, S.T., M.T.

NIP. x

Dr. Alham Fikri Aji, S.T., M.Sc.

NIP. x

DAFTAR ISI

DAFTAR GAMBAR	v
DAFTAR TABEL	vi
DAFTAR KODE	vii
I PENDAHULUAN	1
I.1 Latar Belakang	1
I.2 Rumusan Masalah	3
I.3 Tujuan Penelitian	3
I.4 Batasan Masalah	4
I.5 Metodologi	4
I.5.1 Tahap 1: Investigasi Awal dan Pengumpulan Fakta	5
I.5.2 Tahap 2: Pencarian, Pengelompokan, dan Penapisan Literatur	5
II STUDI LITERATUR	7
II.1 Large Language Model	7
II.1.1 Konsep dan Karakteristik Dasar	7
II.1.2 Representasi Bahasa dan Pemahaman Instruksi	8
II.1.3 Penalaran dan Dinamika Perilaku Model	8
II.1.4 Dimensi Sosial dalam Pemrosesan Bahasa	9
II.2 Persona dalam Interaksi Model Bahasa	10
II.2.1 Definisi dan Ruang Lingkup Persona	10
II.2.2 Persona Eksplisit dan Persona Implisit	10
II.2.3 Peran Persona dalam Interaksi dengan LLM	11
II.3 Pengaruh Persona terhadap Perilaku LLM	12
II.3.1 Pengaruh Persona terhadap Penalaran Model	12
II.3.2 Pengaruh Persona terhadap Gaya dan Struktur Respons	13
II.3.3 Faktor yang Memperkuat Efek Persona	13
II.4 Bias dalam Respons LLM	14
II.4.1 Bentuk-bentuk Bias pada Model Bahasa	14
II.4.2 Konsekuensi Bias terhadap Keluaran Model	15
II.4.3 Kaitannya dengan Variasi Persona	16
II.5 Evaluasi Penalaran dan Benchmark	16
II.5.1 Benchmark Penalaran dan Pengetahuan	17
II.5.2 Benchmark Sosial dan Moral	17

II.5.3	Tantangan Evaluasi Berbasis Persona	18
II.6	Penelitian Terdahulu dan Kesenjangan Penelitian	18
II.6.1	Ringkasan Literatur Terkait	19
II.6.2	Keterbatasan Penelitian Sebelumnya	20
II.6.3	Posisi dan Kontribusi Penelitian Ini	20
II.7	Penulisan Gambar, Tabel, Rumus, dan Kode	21
II.7.1	Gambar	21
II.7.2	Tabel	22
II.7.2.1	Tabel yang Muat dalam Satu Halaman	22
II.7.2.2	Mengimpor Tabel dari Berkas Eksternal	23
II.7.2.3	Tabel yang Sangat Panjang	23
II.7.2.4	Beberapa Contoh Penulisan Rumus atau Persamaan Matematika Menggunakan LaTeX Termasuk Penomorannya	25
II.7.3	Algoritma, Pseudocode, atau Kode	26
II.8	Beberapa Kesalahan Penulisan yang Sering Terjadi	27
II.8.1	Penggunaan Kata "di mana" atau "dimana"	27
II.8.2	Penggunaan Kata "sedangkan" dan "sehingga"	27
II.8.3	Penggunaan Istilah yang Tidak Baku	27
II.8.4	Pemisah Desimal dan Ribuan	28
II.8.5	Daftar atau <i>List</i>	28
II.8.6	Penggunaan Kata "masing-masing" dan "setiap"	28
III	ANALISIS MASALAH	29
III.1	Analisis Kondisi Saat Ini	29
III.2	Analisis Kebutuhan	29
III.2.1	Identifikasi Masalah Pengguna	29
III.2.2	Kebutuhan Fungsional	30
III.2.3	Kebutuhan Nonfungsional	30
III.3	Analisis Pemilihan Solusi	30
III.3.1	Alternatif Solusi	30
III.3.2	Analisis Penentuan Solusi	30
IV	DESAIN KONSEP SOLUSI	32
V	RENCANA SELANJUTNYA	33

DAFTAR GAMBAR

II.1	Contoh gambar jaringan	22
------	----------------------------------	----

DAFTAR TABEL

II.1	Tabel harga bahan pokok	22
II.2	Tabel harga bahan sekunder	23
II.3	Tabel harga bahan tertier	23
II.4	Comprehensive Data Table Example	23
II.4	Comprehensive Data Table Example (lanjutan)	24
II.4	Comprehensive Data Table Example (lanjutan)	25
II.5	Contoh penggunaan kata ”sedangkan” dan ”sehingga”	27

DAFTAR KODE

II.1	Contoh pseudocode	26
II.2	Contoh source code Python	26

BAB I

PENDAHULUAN

I.1 Latar Belakang

Kemajuan dalam pengembangan *large language model* dalam beberapa tahun terakhir telah mengubah cara sistem komputasi memahami, memproses, dan menghasilkan bahasa alami. Model seperti GPT, LLaMA, Mistral, dan Gemini dilatih menggunakan korpus dalam skala masif dan mampu menyelesaikan berbagai tugas mulai dari penalaran numerik hingga interpretasi skenario sosial. Dalam banyak kasus, model menunjukkan kemampuan yang mendekati atau bahkan melampaui performa manusia pada benchmark tertentu. Walaupun demikian, peningkatan kapabilitas ini tidak sepenuhnya diikuti oleh stabilitas perilaku model dalam konteks interaksi dunia nyata.

Salah satu fenomena yang semakin banyak diamati dalam penelitian mutakhir adalah bahwa perilaku *large language model* tidak hanya dipengaruhi oleh isi instruksi, tetapi juga oleh identitas pengguna yang tersirat atau dinyatakan secara eksplisit dalam konteks percakapan. Studi mengenai bias penalaran implisit menunjukkan bahwa perubahan kecil pada deskripsi identitas pengguna dapat menyebabkan variasi signifikan pada hasil penalaran, bahkan untuk tugas yang tidak memiliki aspek sosial eksplisit (Gupta dkk. 2024). Variasi ini mencakup perubahan langkah penalaran, perbedaan tingkat kehati-hatian, hingga munculnya bias tertentu terhadap kelompok sosial.

Selain *user persona* eksplisit yang dituliskan secara langsung dalam instruksi, penelitian menunjukkan bahwa model juga sensitif terhadap *user persona* implisit yang muncul melalui gaya bahasa, framing naratif, struktur pertanyaan, atau atribut linguistik lainnya (Tseng dkk. 2024). Dalam kondisi tersebut, model tidak menerima instruksi tentang identitas pengguna, tetapi tetap membentuk asumsi internal mengenai siapa pengguna dan menyesuaikan respons sesuai asumsi tersebut. Sensitivitas

ini menandakan bahwa model melakukan inferensi identitas pengguna berdasarkan sinyal linguistik yang tampak sepele, yang berimplikasi pada stabilitas penalaran dan keadilan respons.

Penelitian pada bidang pemodelan pengguna menunjukkan bahwa variasi identitas pengguna—seperti usia, latar belakang profesional, afiliasi budaya, atau posisi sosial—dapat memengaruhi keluaran model dalam berbagai dimensi, termasuk penalaran, preferensi jawaban, dan konsistensi respons (Naous, Roziere, dkk. 2025). Hal ini menunjukkan bahwa identitas pengguna, baik eksplisit maupun implisit, berfungsi sebagai variabel laten yang memengaruhi proses generatif model. Dengan demikian, analisis terhadap *user persona* menjadi penting tidak hanya untuk memahami perilaku model, tetapi juga untuk mengidentifikasi potensi bias dan ketidakstabilan yang muncul dalam interaksi manusia–AI.

Walaupun berbagai studi sebelumnya memberikan indikasi bahwa identitas pengguna memengaruhi perilaku model, penelitian yang ada masih memiliki batasan. Mayoritas studi hanya mengevaluasi satu atau dua model, cakupan persona yang terbatas, atau jenis tugas yang sempit. Selain itu, tidak banyak studi yang secara sistematis membandingkan efek *user persona* eksplisit dan implisit pada berbagai model dan berbagai jenis tugas penalaran dalam satu kerangka eksperimen yang konsisten. Belum tersedia pula pendekatan evaluasi yang secara terpadu menguji sensitivitas model terhadap variasi identitas pengguna di berbagai kondisi tugas, baik numerik, logis, faktual, sosial, maupun moral.

Kekosongan penelitian ini penting untuk dijawab, mengingat model bahasa semakin banyak digunakan pada skenario yang sensitif terhadap identitas pengguna, seperti layanan kesehatan, pendidikan, konseling, sistem rekomendasi, dan interaksi berbasis nilai. Ketidakstabilan respons akibat identitas pengguna berpotensi menimbulkan bias, mengurangi keandalan model, dan menghasilkan ketidaksetaraan dalam pengalaman pengguna. Oleh karena itu, diperlukan pendekatan evaluasi yang lebih komprehensif untuk memahami bagaimana *user persona* eksplisit dan implisit memengaruhi penalaran, perilaku keluaran, dan kecenderungan *human bias* pada berbagai *large language model*.

Berdasarkan urgensi tersebut, penelitian ini disusun untuk melakukan evaluasi empiris terhadap pengaruh *user persona* eksplisit dan *user persona* implisit melalui eksperimen terstruktur pada berbagai model dan berbagai jenis tugas. Penelitian ini diharapkan memberikan pemahaman yang lebih mendalam mengenai sensitivitas model terhadap identitas pengguna serta implikasinya terhadap penalaran, bias, dan

keandalan model dalam aplikasi dunia nyata.

I.2 Rumusan Masalah

Rumusan masalah berikut disusun berdasarkan kebutuhan untuk memahami bagaimana *user persona* memengaruhi perilaku dan penalaran model bahasa. Penelitian sebelumnya menunjukkan bahwa identitas pengguna, baik yang diberikan secara eksplisit maupun implisit, dapat memengaruhi penalaran, kualitas keluaran, dan kecenderungan bias model (Gupta dkk. 2024; Tseng dkk. 2024; Naous, Roziere, dkk. 2025). Namun, cakupan penelitian terdahulu masih terbatas pada sedikit model, sedikit persona, dan variasi tugas yang sempit.

Berdasarkan kondisi tersebut, rumusan masalah penelitian ini adalah sebagai berikut.

1. Bagaimana pengaruh *user persona* eksplisit dan *user persona* implisit terhadap performa penalaran pada berbagai jenis tugas pada sejumlah *large language model*.
2. Bagaimana kedua jenis *user persona* tersebut memengaruhi perilaku keluaran model pada skenario interaksi yang berbeda.
3. Bagaimana pola *human bias* muncul dan berubah sebagai akibat variasi *user persona*.
4. Sejauh mana sensitivitas terhadap *user persona* berbeda pada berbagai *large language model*, serta model mana yang menunjukkan tingkat *robustness* yang lebih tinggi terhadap variasi tersebut.

I.3 Tujuan Penelitian

Tujuan penelitian ditetapkan untuk menjawab permasalahan yang telah dirumuskan. Penelitian ini diarahkan untuk menghasilkan pemahaman yang lebih komprehensif mengenai pengaruh *user persona* terhadap perilaku model bahasa dalam tugas penalaran dan skenario percakapan. Secara khusus, penelitian ini bertujuan untuk:

1. Menganalisis pengaruh *user persona* eksplisit dan *user persona* implisit terhadap performa penalaran pada sejumlah *large language model*.
2. Mengidentifikasi perubahan perilaku keluaran model yang diinduksi oleh variasi *user persona* pada berbagai konteks.
3. Menganalisis pola *human bias* yang muncul akibat variasi *user persona*.
4. Menyusun perbandingan sensitivitas dan *robustness* berbagai model terhadap variasi *user persona*.

5. Mengembangkan rancangan *evaluation pipeline* yang memungkinkan pelaksanaan eksperimen *multi model* dan *multi persona* secara terotomatisasi.

I.4 Batasan Masalah

Batasan masalah ditetapkan agar ruang lingkup penelitian terkelola dan selaras dengan tujuan penelitian. Penelitian ini tidak bertujuan mengevaluasi seluruh aspek perilaku model bahasa, tetapi fokus pada pengaruh *user persona*. Batasan penelitian ini adalah sebagai berikut.

1. Penelitian hanya menganalisis dua jenis *user persona*, yaitu *user persona* eksplisit dan *user persona* implisit. Penelitian tidak mencakup *role-playing persona* yang memberikan identitas kepada model maupun mekanisme *personalization* berbasis histori pengguna.
2. Pengujian terbatas pada model bahasa berbasis teks yang dapat diakses melalui API. Model multimodal, model yang memerlukan *fine-tuning*, atau model yang memerlukan pelatihan ulang tidak termasuk dalam cakupan penelitian.
3. Evaluasi dibatasi pada tugas berbasis teks, termasuk penalaran numerik, penalaran logis, pertanyaan pengetahuan umum, skenario sosial, dan skenario moral. Tugas vision-language atau *speech* tidak dibahas.
4. Penilaian kualitas keluaran dilakukan melalui evaluasi terotomatisasi dan analisis komparatif. Penilaian berbasis partisipan manusia tidak dilakukan.
5. Penelitian menggunakan *evaluation pipeline* berbasis eksekusi prompt tanpa melakukan modifikasi pada parameter internal model.
6. Analisis bias terbatas pada *human bias* yang muncul sebagai akibat variasi *user persona*, dan tidak mencakup bias makro yang bersumber dari data pelatihan model.

I.5 Metodologi

Metodologi pada tahap penyusunan proposal ini disusun untuk memastikan bahwa proses perumusan masalah, penentuan ruang lingkup penelitian, dan penyusunan kerangka teoretis dilakukan secara sistematis. Metodologi ini tidak mencakup tahapan implementasi eksperimen, yang akan dijabarkan pada Bab III, melainkan berfokus pada kegiatan awal yang diperlukan untuk menghasilkan proposal penelitian yang terarah dan berbasis kajian ilmiah.

I.5.1 Tahap 1: Investigasi Awal dan Pengumpulan Fakta

Tahap awal dilakukan untuk memahami konteks permasalahan dan mengidentifikasi isu ilmiah yang relevan dengan topik penelitian. Langkah yang dilakukan meliputi:

1. Mengidentifikasi fenomena sensitivitas *large language model* terhadap identitas pengguna berdasarkan contoh kasus, laporan empiris, dan temuan penelitian sebelumnya.
2. Meninjau keluaran awal beberapa model bahasa melalui eksplorasi terbatas untuk mengamati indikasi pengaruh *user persona* eksplisit dan *user persona* implisit terhadap penalaran dan gaya respons.
3. Menyimpulkan pola permasalahan yang muncul untuk kemudian dirumuskan sebagai pokok masalah penelitian.

I.5.2 Tahap 2: Pencarian, Pengelompokan, dan Penapisan Literatur

Tahap ini dilakukan untuk memperoleh landasan ilmiah yang kuat dalam menyusun kerangka teoretis dan menentukan arah penelitian. Kegiatan yang dilakukan mencakup:

1. Melakukan pencarian literatur menggunakan mesin pencarian akademik seperti Google Scholar, Semantic Scholar, arXiv, dan ACL Anthology dengan kata kunci antara lain *user persona*, *implicit persona*, *identity-conditioned prompting*, *LLM sensitivity*, *reasoning evaluation*, dan *bias in LLM*.
2. Menyeleksi publikasi yang relevan, termasuk penelitian mengenai pengaruh identitas pengguna terhadap keluaran model bahasa, teori penalaran pada model bahasa, evaluasi berbasis prompt, dan bias implisit.
3. Mengelompokkan literatur ke dalam kategori konseptual, yaitu: (a) konsep dasar *large language model*, (b) teori dan klasifikasi *persona* eksplisit dan implisit, (c) penelitian terdahulu mengenai identitas pengguna dan pengaruhnya terhadap keluaran model, (d) metode evaluasi penalaran dan analisis bias.
4. Menganalisis dan merangkum kontribusi, metodologi, serta keterbatasan setiap publikasi yang terpilih untuk memastikan bahwa kerangka teoretis proposal didasarkan pada referensi yang valid dan mutakhir.
5. Mendokumentasikan seluruh proses penelusuran literatur, termasuk daftar kata kunci, sumber pencarian, dan kriteria penapisan yang digunakan. Dokumentasi tambahan, seperti rekaman proses eksplorasi awal atau catatan observasi, akan dicantumkan pada bagian lampiran.

Tahap-tahap tersebut menghasilkan landasan konseptual dan rumusan permasalahan yang digunakan dalam penyusunan proposal tugas akhir. Hasil kajian literatur secara

rinci akan disajikan pada Bab II Studi Literatur.

BAB II

STUDI LITERATUR

Bab ini membahas konsep dan penelitian terdahulu yang menjadi landasan bagi analisis pengaruh *user persona* terhadap perilaku *large language model*. Pembahasan disusun secara bertahap, dimulai dari uraian mengenai model bahasa modern, mekanisme pemrosesan instruksi, konsep dasar persona, serta temuan empiris mengenai sensitivitas model terhadap identitas pengguna. Selain itu, bab ini meninjau isu bias dan metode evaluasi penalaran yang relevan bagi perancangan penelitian ini.

II.1 Large Language Model

II.1.1 Konsep dan Karakteristik Dasar

Large language model (LLM) merupakan model generatif berbasis arsitektur transformator yang dilatih menggunakan data dalam skala sangat besar. Model ini mempelajari pola bahasa melalui hubungan antartoken, sehingga mampu membangun representasi yang mencakup makna, hubungan semantik, serta isyarat pragmatik yang muncul dalam teks. Dengan skala pelatihan yang luas, LLM dapat digunakan pada berbagai tugas tanpa memerlukan penyesuaian khusus untuk setiap tugas.

Secara konseptual, LLM bekerja dengan memprediksi token berikutnya berdasarkan konteks sebelumnya. Namun, proses prediksi ini tidak sekadar berbasis frekuensi kata, melainkan menggunakan representasi kontekstual yang memungkinkan model memahami instruksi, gaya penulisan, maupun kecenderungan komunikasi. Model seperti GPT, LLaMA, Mistral, dan Gemini mengadopsi pendekatan ini dan menunjukkan kemampuan generalisasi yang kuat terhadap tugas bahasa yang kompleks.

Karakteristik utama LLM antara lain fleksibilitas dalam mengikuti instruksi, kemampuan menyusun penalaran, serta penyesuaian terhadap pola komunikasi pengguna. Kemampuan ini muncul dari kombinasi arsitektur dasar transformator, skala

parameter yang besar, dan keragaman data pelatihan. Karena model tidak dibuat untuk satu domain tertentu, tetapi dilatih pada data lintas konteks, gaya, dan situasi, LLM dapat mengadaptasi perilaku komunikasinya berdasarkan variasi kecil dalam instruksi.

II.1.2 Representasi Bahasa dan Pemahaman Instruksi

LLM memproses teks melalui beberapa tahapan representasi internal. Teks diuraikan menjadi token, kemudian dipetakan ke dalam ruang representasi berdimensi tinggi melalui *embedding*. Representasi awal ini kemudian diperkaya melalui lapisan-lapisan transformator yang memanfaatkan mekanisme perhatian untuk menentukan hubungan antar token dalam konteks yang lebih luas. Hasilnya adalah representasi kontekstual yang mencerminkan interpretasi model terhadap instruksi atau percakapan.

Representasi ini tidak bersifat statis. Makna sebuah token dapat berubah bergantung pada cara pengguna menyampaikan instruksi. Perbedaan gaya penulisan, urutan informasi, atau tingkat formalitas dapat menghasilkan representasi internal yang berbeda, sehingga memunculkan respons yang berbeda pula. Penelitian Zhou et al. (Zhou dkk. 2023) menunjukkan bahwa perubahan kecil dalam framing, seperti perbedaan nada atau cara bertanya, dapat menggeser perhatian model dan mengubah struktur jawaban yang dihasilkan.

Sebagai ilustrasi, perbedaan instruksi berikut sering kali menghasilkan respons yang berbeda meskipun inti pertanyaannya sama:

- “Jelaskan secara singkat apa itu regularisasi.”
- “Saya sedang menulis laporan akademik. Bisakah Anda menjelaskan secara formal apa yang dimaksud dengan regularisasi?”

Instruksi kedua biasanya memicu model untuk memberikan penjelasan yang lebih panjang, lebih berhati-hati, dan lebih formal. Perbedaan ini mencerminkan bagaimana representasi instruksi terbentuk berdasarkan konteks linguistik dan pragmatik.

II.1.3 Penalaran dan Dinamika Perilaku Model

Selain pemahaman instruksi, LLM juga menunjukkan kemampuan melakukan penalaran. Model dapat menyelesaikan soal penalaran numerik sederhana, menjawab pertanyaan berbasis pengetahuan umum, hingga memberikan penilaian terhadap skenario sosial atau moral. Namun, kemampuan ini tidak sepenuhnya stabil. Turpin et al. (Turpin dkk. 2023) menemukan bahwa penalaran yang dihasilkan model

dapat berubah hanya karena variasi kecil pada bentuk instruksi, walaupun substansi tugas tetap sama.

Hal ini terjadi karena model tidak melakukan penalaran melalui prosedur logis eksplisit, tetapi melalui dinamika representasi internal yang sensitif terhadap konteks. Sebuah instruksi yang lebih panjang atau lebih formal dapat memicu struktur penalaran yang lebih sistematis, sementara instruksi yang lebih langsung dapat menghasilkan jawaban tanpa uraian langkah-langkah penalaran yang jelas. Perubahan ini memperlihatkan bahwa struktur penalaran yang muncul merupakan fungsi dari konteks interaksi, bukan semata-mata fungsi dari logika masalah yang diberikan.

Ketidakstabilan ini penting untuk dipahami karena berhubungan langsung dengan penelitian mengenai *user persona*. Jika perubahan kecil pada instruksi dapat mengubah penalaran, maka variasi identitas pengguna yang tersirat dalam tulisan juga berpotensi memicu perubahan serupa.

II.1.4 Dimensi Sosial dalam Pemrosesan Bahasa

Model bahasa modern tidak hanya mempelajari struktur dan makna bahasa, tetapi juga pola interaksi sosial yang tercermin dalam data pelatihan. Weidinger et al. (Weidinger dkk. 2021) menunjukkan bahwa LLM dapat menginternalisasi norma sosial, stereotip, serta pola komunikasi yang umum digunakan manusia. Dalam banyak kasus, gaya bahasa tertentu diinterpretasikan sebagai sinyal sosial mengenai siapa pengguna tersebut, misalnya usia, latar profesional, atau tingkat pendidikan.

Ketika instruksi ditulis dengan gaya santai, model sering kali memberikan respons yang lebih ringkas atau lebih langsung. Sebaliknya, ketika instruksi ditulis dengan gaya formal, respons yang dihasilkan cenderung lebih berhati-hati dan mengikuti struktur penjelasan akademis. Perbedaan respons ini bukan sekadar akibat gaya penulisan, tetapi akibat inferensi sosial yang dilakukan model berdasarkan pola komunikasi dalam data pelatihan.

Fenomena ini menunjukkan bahwa pemrosesan bahasa oleh LLM memiliki dimensi sosial yang signifikan. Instruksi diperlakukan bukan hanya sebagai teks, tetapi sebagai bentuk interaksi manusia yang membawa sinyal identitas. Sensitivitas terhadap sinyal ini merupakan salah satu alasan mengapa *user persona* dapat memengaruhi penalaran, struktur respons, maupun kecenderungan bias dalam keluaran model.

II.2 Persona dalam Interaksi Model Bahasa

II.2.1 Definisi dan Ruang Lingkup Persona

Dalam kajian sistem bahasa alami, *persona* merujuk pada serangkaian atribut yang digunakan untuk menggambarkan identitas atau karakteristik pengguna. Atribut tersebut dapat berupa informasi sosial, demografis, profesional, atau gaya komunikasi yang merepresentasikan cara seseorang berinteraksi dalam percakapan. Persona berfungsi sebagai konteks tambahan yang dapat memengaruhi bagaimana sebuah sistem dialog memahami maksud pengguna dan membentuk respons.

Dalam konteks *large language model*, persona tidak hanya dipandang sebagai label identitas, tetapi juga sebagai bagian dari sinyal yang terkandung dalam bahasa. Karena model belajar dari data pelatihan yang mencerminkan cara manusia berkomunikasi, model juga mempelajari keterkaitan antara gaya bahasa dan identitas sosial. Dengan demikian, persona tidak hanya bekerja sebagai informasi eksplisit, tetapi dapat tersirat melalui variasi linguistik seperti pilihan kata, nada, struktur kalimat, atau keformalan tulisan.

Ruang lingkup persona dalam sistem bahasa mencakup berbagai kategori identitas, seperti gender, usia, minat, latar profesional, afiliasi budaya, ataupun preferensi komunikasi. Representasi persona tersebut tidak selalu hadir dalam bentuk pernyataan langsung, tetapi sering kali dinyatakan melalui konteks linguistik yang halus tanpa deklarasi eksplisit mengenai siapa pengguna tersebut.

II.2.2 Persona Eksplisit dan Persona Implisit

Fenomena persona dalam interaksi dengan model bahasa dapat dibagi menjadi dua bentuk utama, yaitu persona eksplisit dan persona implisit. Keduanya memberikan sinyal identitas, tetapi melalui mekanisme dan intensitas yang berbeda.

Persona eksplisit muncul ketika identitas pengguna dinyatakan secara langsung dalam instruksi atau konteks percakapan. Contohnya adalah ketika pengguna menuliskan “Saya adalah mahasiswa teknik informatika” atau “Sebagai seorang dokter, saya ingin memahami...”. Ungkapan seperti ini memberikan sinyal yang jelas kepada model mengenai latar pengguna, sehingga model dapat menyesuaikan struktur respons agar lebih sesuai dengan karakteristik tersebut. Gupta et al. (Gupta dkk. 2024) menunjukkan bahwa penugasan persona eksplisit semacam ini dapat mengubah hasil penalaran model, meskipun tugas yang diberikan tidak berkaitan dengan identitas sosial pengguna. Perubahan respons tidak hanya menyangkut gaya bahasa,

tetapi juga dapat memengaruhi kesimpulan logis yang diberikan model.

Sebaliknya, persona implisit muncul ketika identitas pengguna tidak dinyatakan secara langsung, tetapi disimpulkan oleh model berdasarkan isyarat linguistik. Penelitian Tseng et al. (Tseng dkk. 2024) menunjukkan bahwa model memiliki kecenderungan melakukan inferensi identitas pengguna dari gaya penulisan, struktur kalimat, pilihan kata, atau tingkat formalitas. Fenomena ini dapat terjadi meskipun pengguna tidak bermaksud menyampaikan identitas tertentu. Sebagai contoh, gaya penulisan formal dengan istilah akademis sering diasosiasikan dengan latar pendidikan tertentu, sedangkan gaya penulisan santai dapat diasosiasikan dengan kategori usia atau tingkat kedekatan sosial.

Inferensi identitas tersebut bukan hasil dari aturan yang ditetapkan secara eksplisit dalam model, tetapi merupakan konsekuensi dari pola komunikasi manusia yang terserap selama proses pelatihan. Model mempelajari bahwa gaya bahasa tertentu sering muncul bersama atribut sosial tertentu, sehingga ketika gaya tersebut muncul dalam instruksi, model cenderung mengaktifkan pola respons yang sesuai dengan kategori identitas yang diasosiasikan. Fenomena ini menjadi dasar penting bagi studi mengenai pengaruh persona implisit terhadap perilaku dan penalaran model.

II.2.3 Peran Persona dalam Interaksi dengan LLM

Persona, baik eksplisit maupun implisit, berperan sebagai sinyal kontekstual yang memengaruhi interpretasi dan respons model bahasa. Ketika identitas pengguna muncul dalam bentuk atribut sosial atau gaya komunikasi tertentu, model akan memperlakukannya sebagai bagian dari konteks yang relevan. Konteks ini kemudian membentuk representasi internal yang memengaruhi bagaimana model memahami pertanyaan, menafsirkan maksud, dan menyusun jawaban.

Peran persona dalam interaksi ini dapat dilihat dari dua dimensi utama. Pertama, persona dapat memengaruhi aspek linguistik respons, seperti pilihan kata, tingkat formalitas, pola argumentasi, atau struktur penjelasan. Model cenderung menyesuaikan respons agar selaras dengan gaya komunikasi yang diasosiasikan dengan persona tertentu. Kedua, persona dapat memengaruhi penalaran model melalui apa yang disebut sebagai *reasoning shift*, yaitu perubahan struktur penalaran yang terjadi akibat variasi identitas pengguna meskipun substansi tugas tetap sama.

Sebagai ilustrasi, suatu pertanyaan logika sederhana yang diajukan oleh pengguna dengan persona profesional tertentu dapat memicu model untuk memberikan res-

pons yang lebih sistematis atau lebih berhati-hati. Sebaliknya, pertanyaan yang diajukan dengan gaya informal dapat menghasilkan respons yang lebih ringkas dengan struktur penalaran minimal. Perubahan ini menunjukkan bahwa persona berfungsi sebagai variabel kondisi yang membentuk dinamika interaksi antara pengguna dan model.

II.3 Pengaruh Persona terhadap Perilaku LLM

Pembahasan mengenai persona tidak berhenti pada bagaimana identitas pengguna direpresentasikan dalam instruksi, tetapi juga mencakup bagaimana identitas tersebut memengaruhi perilaku model bahasa ketika menghasilkan respons. Berbagai penelitian menunjukkan bahwa persona berperan sebagai konteks tambahan yang secara halus membentuk cara model memahami pertanyaan, menimbang informasi, dan menyusun jawaban. Dengan demikian, persona tidak sekadar menjadi atribut linguistik, tetapi menjadi bagian dari dinamika interaksi yang memengaruhi proses penalaran dan karakter keluaran model.

II.3.1 Pengaruh Persona terhadap Penalaran Model

Penalaran merupakan salah satu kemampuan utama yang ditonjolkan oleh model bahasa modern. Namun, sejumlah studi menemukan bahwa penalaran tersebut tidak selalu stabil dan dapat berubah bergantung pada konteks identitas pengguna. Gupta et al. (Gupta dkk. 2024) menunjukkan bahwa ketika sebuah persona eksplisit disisipkan ke dalam instruksi, model dapat menghasilkan struktur penalaran yang berbeda meskipun tugas yang diberikan tetap sama. Perubahan tersebut terlihat pada pemilihan langkah-langkah argumentatif, urutan penjelasan, atau tingkat kehati-hatian dalam menarik kesimpulan.

Dalam konteks persona implisit, perubahan penalaran muncul melalui mekanisme yang lebih halus. Gaya penulisan pengguna, seperti tingkat formalitas, panjang kalimat, atau pilihan kosakata, dapat diinterpretasikan sebagai sinyal identitas yang memengaruhi cara model membangun penalaran. Misalnya, instruksi yang disampaikan dengan gaya akademis sering kali mendorong model untuk memberikan penjelasan yang lebih sistematis dan rinci. Sebaliknya, instruksi yang ditulis dengan gaya santai dapat menghasilkan penalaran yang lebih ringkas atau langsung.

Temuan-temuan ini sejalan dengan penelitian mengenai ketidakstabilan penalaran yang dilakukan oleh Turpin et al. (Turpin dkk. 2023). Dalam studi tersebut, perubahan kecil pada struktur instruksi terbukti memengaruhi urutan *chain-of-thought*

yang dihasilkan model. Karena persona bekerja sebagai bagian dari konteks instruksi, variasi identitas pengguna berpotensi menimbulkan pergeseran pola berpikir yang muncul dalam respons model.

Pengaruh persona terhadap penalaran tampak pada berbagai kategori tugas, mulai dari penalaran numerik hingga pertimbangan moral. Pada tugas numerik, persona tertentu dapat mendorong model untuk memberikan uraian langkah yang lebih panjang atau lebih hati-hati. Pada tugas logika, persona dapat memengaruhi cara model menyusun argumen. Sementara itu, pada tugas sosial atau moral, persona dapat mengarahkan model untuk menekankan nilai-nilai tertentu atau memilih perspektif yang lebih dekat dengan identitas pengguna yang diasumsikan.

II.3.2 Pengaruh Persona terhadap Gaya dan Struktur Respons

Selain penalaran, persona juga memengaruhi aspek gaya dan struktur respons. Model bahasa modern tidak hanya menghasilkan jawaban berdasarkan isi pertanyaan, tetapi juga menyesuaikan cara penyampaiannya agar selaras dengan identitas pengguna yang terdeteksi. Temuan Tseng et al. (Tseng dkk. 2024) menunjukkan bahwa model dapat meniru gaya bahasa yang diasosiasikan dengan persona tertentu, bahkan ketika identitas tersebut tidak dinyatakan secara eksplisit.

Perubahan yang muncul dapat berupa pemilihan kosakata, panjang penjelasan, tingkat formalitas, atau nada yang digunakan dalam respons. Apabila model mengaitkan pengguna dengan latar profesional tertentu, respons yang dihasilkan sering kali lebih teknis atau lebih terstruktur. Sebaliknya, apabila gaya penulisan pengguna menunjukkan kedekatan sosial atau informalitas, respons yang muncul cenderung lebih ringkas atau lebih langsung.

Dalam beberapa kasus, persona tertentu juga dapat memicu model untuk bersikap lebih berhati-hati, terutama pada topik-topik yang sensitif secara sosial. Fenomena ini muncul karena model mempelajari pola komunikasi manusia dalam data pelatihan dan mengaitkan gaya bahasa dengan norma sosial yang berlaku pada kelompok tertentu. Dengan demikian, perbedaan gaya respons bukan sekadar variasi permukaan, tetapi merupakan hasil dari proses interpretasi sosial yang dilakukan model.

II.3.3 Faktor yang Memperkuat Efek Persona

Variasi respons akibat persona diperkuat oleh sejumlah faktor yang berkaitan dengan konteks interaksi. Salah satu faktor tersebut adalah framing instruksi. Ketika persona disampaikan secara konsisten, baik melalui deskripsi eksplisit maupun ga-

ya penulisan yang stabil, representasi identitas pengguna menjadi lebih kuat dalam interpretasi model. Hal ini membuat model lebih cenderung mempertahankan pola respons tertentu sepanjang percakapan.

Selain itu, jenis tugas yang diberikan turut memengaruhi seberapa besar dampak persona terhadap respons model. Tugas yang bersifat terbuka, seperti pertanyaan moral atau skenario sosial, memberikan ruang interpretasi yang lebih luas sehingga sinyal identitas lebih mudah memengaruhi pola jawaban. Sebaliknya, tugas-tugas yang memiliki jawaban pasti atau struktur penyelesaian yang ketat cenderung menunjukkan pengaruh persona yang lebih kecil.

Skala model dan metode penyesuaian instruksi juga memainkan peran penting. Model yang dilatih dengan data percakapan dalam jumlah besar cenderung lebih sensitif terhadap variasi gaya linguistik. Sementara itu, model dengan kapasitas lebih kecil dapat menunjukkan respons yang kurang konsisten karena representasi sosial yang terbatas.

Secara keseluruhan, efek persona merupakan hasil interaksi antara konteks linguistik, representasi sosial, dan mekanisme penyesuaian model. Faktor-faktor ini bekerja bersamaan dan membentuk variasi respons yang menggambarkan bagaimana model bahasa menafsirkan identitas pengguna dalam proses menghasilkan jawaban.

II.4 Bias dalam Respons LLM

Pembahasan mengenai bias dalam *large language model* berangkat dari kenyataan bahwa model bahasa belajar dari pola-pola yang muncul dalam data pelatihan. Data tersebut bukan hanya berisi informasi faktual, tetapi juga memuat kecenderungan sosial yang terbentuk secara historis. Ketika model mempelajari pola bahasa dari data tersebut, model tidak hanya menyerap struktur linguistik, tetapi juga asumsi-asumsi sosial yang secara tidak sengaja dapat tercermin dalam respons yang dihasilkan. Dalam konteks penelitian ini, bias menjadi penting karena persona—baik eksplisit maupun implisit—dapat memperkuat atau menggeser pola bias yang dimiliki model.

II.4.1 Bentuk-bentuk Bias pada Model Bahasa

Bias dalam model bahasa dapat muncul dalam beragam bentuk. Salah satu bentuk yang sering menjadi perhatian adalah bias representasional, yaitu kecenderungan model menggambarkan suatu kelompok sosial secara tidak seimbang. Weidinger et

al. (Weidinger dkk. 2021) menunjukkan bahwa stereotip yang sering muncul dalam teks internet dapat terinternalisasi dalam model. Misalnya, model dapat mengaitkan profesi tertentu dengan gender tertentu atau menempatkan kelompok sosial tertentu dalam peran tertentu, meskipun konteks yang diberikan sebenarnya netral.

Selain bias representasional, terdapat bias inferensial, yaitu kecenderungan model mengambil kesimpulan berdasarkan isyarat yang tidak relevan. Bentuk bias ini biasanya muncul ketika model meniru pola asosiasi dari data tanpa memahami konteks sebenarnya. Sebagai contoh, ketika diminta mendeskripsikan seseorang dalam skenario imajinatif, model dapat mengisi detail yang tidak disebutkan hanya karena mengikuti pola umum yang sering ditemui dalam data pelatihan.

Bias tersebut tidak hanya muncul pada isi jawaban, tetapi juga dalam cara model menyusun uraian penjelasan. Pada topik-topik moral atau sosial, bias dapat terlihat dari pilihan nilai atau asumsi yang digunakan dalam penalaran. Hal ini memperlihatkan bahwa bias dalam model bahasa bersifat berlapis: ia dapat memengaruhi kosakata, struktur kalimat, hingga cara model melakukan evaluasi terhadap suatu situasi.

II.4.2 Konsekuensi Bias terhadap Keluaran Model

Bias yang muncul dalam model bahasa membawa sejumlah konsekuensi terhadap keluaran yang diberikan kepada pengguna. Salah satu konsekuensi yang paling sering dibahas adalah risiko misinformasi. Ketika model memberikan jawaban yang terdengar meyakinkan tetapi sebenarnya bias atau tidak akurat, pengguna yang tidak memiliki pengetahuan memadai dapat menerima informasi tersebut sebagai kebenaran.

Konsekuensi lainnya berkaitan dengan ketidakmerataan kualitas respons. Jika model menyesuaikan gaya penjelasan berdasarkan persona tertentu, kelompok pengguna yang berbeda dapat menerima penjelasan dengan tingkat kedalaman atau kehati-hatian yang tidak sama. Walaupun model tidak memiliki niat atau tujuan tertentu, perbedaan kualitas informasi ini dapat mempengaruhi proses pemahaman pengguna terhadap suatu topik.

Di sisi lain, bias juga dapat memperkuat stereotip sosial. Ketika model berulang kali memberikan deskripsi atau penilaian yang sejalan dengan stereotip tertentu, model secara tidak langsung ikut berpartisipasi dalam memperkuat persepsi sosial yang tidak akurat. Penguatan stereotip ini dapat terjadi secara halus, misalnya

melalui pilihan kosakata yang cenderung bernuansa tertentu atau struktur argumen yang mengarah pada penilaian yang bias.

II.4.3 Kaitannya dengan Variasi Persona

Persona, sebagai sinyal identitas pengguna, dapat memperkuat atau menggeser munculnya bias dalam respons model. Pada persona eksplisit, bias dapat timbul ketika model mengaitkan identitas pengguna dengan pola stereotip dalam data pelatihan. Misalnya, pernyataan seperti “Saya seorang guru” atau “Saya berasal dari profesi X” dapat memicu model memberikan respons yang mengikuti pola tertentu yang sering dikaitkan dengan profesi tersebut.

Pada persona implisit, bias muncul dengan cara yang lebih halus. Karena model sangat peka terhadap gaya penulisan, pilihan kata atau tingkat formalitas dapat dianggap sebagai indikator identitas sosial pengguna. Jika model mengaitkan gaya komunikasi tertentu dengan kelompok sosial tertentu, respons yang dihasilkan dapat mencerminkan bias yang dimiliki model terhadap kelompok tersebut.

Fenomena ini menjadi lebih terlihat pada tugas-tugas yang bersifat terbuka, seperti pertanyaan moral, skenario etika, atau pertimbangan sosial. Pada jenis tugas tersebut, respons model sangat dipengaruhi oleh konteks dan cara model melakukan inferensi sosial. Ketika persona menjadi bagian dari konteks, respons yang dihasilkan dapat menunjukkan pergeseran nilai, perhatian, atau prioritas tertentu. Kondisi inilah yang membuat analisis persona dalam penelitian ini menjadi penting: persona tidak hanya memengaruhi gaya bahasa atau cara penalaran, tetapi juga membuka ruang bagi bias untuk muncul atau berubah.

II.5 Evaluasi Penalaran dan Benchmark

Evaluasi terhadap *large language model* tidak hanya dilakukan dengan melihat kemampuan model menghasilkan teks, tetapi juga melalui serangkaian tugas terstruktur yang dirancang untuk mengukur kemampuan penalaran, pemahaman konteks, serta kemampuan model menyelesaikan masalah secara konsisten. Benchmark menjadi alat penting dalam penelitian karena memberikan gambaran yang lebih objektif mengenai bagaimana model berperilaku di berbagai situasi dan tingkat kesulitan. Dalam konteks penelitian ini, benchmark yang digunakan tidak hanya berfungsi untuk menilai performa penalaran, tetapi juga untuk melihat bagaimana persona dapat memengaruhi keluaran model pada berbagai jenis tugas.

II.5.1 Benchmark Penalaran dan Pengetahuan

Sejumlah benchmark telah dikembangkan untuk mengukur kemampuan penalaran model bahasa. Salah satu yang paling dikenal adalah GSM8K, sebuah kumpulan soal matematika tingkat sekolah dasar yang dirancang untuk menguji penalaran numerik dan kemampuan model menyusun langkah-langkah penyelesaian secara terstruktur. Meskipun soalnya sederhana bagi manusia, benchmark ini cukup menantang bagi model bahasa karena mengharuskan model memahami konteks, menerapkan logika dasar, dan menjaga konsistensi antara uraian langkah dan jawaban akhir.

Selain GSM8K, benchmark lain seperti MMLU digunakan untuk menguji kemampuan model pada pertanyaan lintas domain, mulai dari sains hingga ilmu sosial. MMLU menekankan kapasitas model dalam memahami pengetahuan faktual dan menerapkannya dalam konteks yang tepat. Benchmark ini memberikan gambaran mengenai seberapa baik model dapat menjawab pertanyaan yang membutuhkan pemahaman konsep dan penalaran tingkat menengah.

Benchmark semacam ini penting karena menampilkan kemampuan dasar model tanpa dipengaruhi oleh gaya interaksi yang terlalu terbuka. Dengan kata lain, benchmark berbasis pengetahuan atau logika dasar memberikan titik awal yang netral sebelum mempertimbangkan bagaimana persona dapat menggeser atau memengaruhi jawaban model.

II.5.2 Benchmark Sosial dan Moral

Di samping penalaran numerik dan faktual, kemampuan model untuk memahami situasi sosial dan moral juga menjadi perhatian dalam penelitian. Benchmark seperti SocialIQA digunakan untuk mengukur kemampuan model memahami skenario sosial sederhana, misalnya bagaimana seseorang mungkin merespons suatu tindakan atau apa motivasi yang mungkin dimiliki dalam konteks tertentu. Benchmark ini menekankan bagaimana model menginternalisasi pola interaksi antarindividu berdasarkan data pelatihan.

Selain SocialIQA, terdapat pula tugas-tugas moral yang dirancang untuk melihat bagaimana model memberikan penilaian terhadap situasi etis. Tugas semacam ini tidak memiliki jawaban pasti, sehingga respons model sangat dipengaruhi oleh nilai, norma, atau pola argumentasi yang diserap selama pelatihan. Dalam konteks penelitian persona, tugas moral menjadi menarik karena persona dapat menggeser sudut pandang moral yang diambil model, misalnya apakah model menjadi lebih

berhati-hati, lebih permissive, atau lebih normatif.

Benchmark sosial dan moral ini penting untuk menganalisis bagaimana persona bekerja pada situasi yang tidak memiliki jawaban tunggal dan mengharuskan model melakukan interpretasi berdasarkan konteks sosial.

II.5.3 Tantangan Evaluasi Berbasis Persona

Meskipun benchmark merupakan alat penting dalam evaluasi model, penggunaan benchmark dalam penelitian persona memiliki tantangan tersendiri. Salah satu tantangan utama adalah konsistensi. Karena persona dapat memengaruhi gaya penalaran dan respons model, evaluasi harus dilakukan dengan cara yang memastikan bahwa perubahan yang muncul benar-benar disebabkan oleh persona, bukan oleh variasi lain dalam instruksi atau struktur prompt.

Tantangan berikutnya adalah sensitivitas model terhadap framing. Perubahan kecil pada instruksi, bahkan ketika persona tidak berubah, dapat menghasilkan respons yang berbeda. Hal ini membuat evaluasi berbasis persona memerlukan desain eksperimen yang hati-hati agar pengaruh persona dapat dipisahkan dari pengaruh variasi linguistik.

Selain itu, model bahasa cenderung mengalami *drift* atau perubahan perilaku kecil antarevaluasi, terutama jika evaluasi tidak terotomatisasi dengan baik. Hal ini dapat memengaruhi replikasi hasil dan interpretasi terhadap pengaruh persona. Penggunaan *pipeline* evaluasi yang terstandardisasi dapat membantu mengurangi variasi ini dengan memastikan bahwa setiap model menerima struktur instruksi yang konsisten.

Secara keseluruhan, benchmark memberikan fondasi penting untuk memahami perilaku model dalam berbagai skenario. Namun, dalam konteks penelitian persona, benchmark tidak hanya berfungsi sebagai alat ukur kemampuan, tetapi juga sebagai sarana untuk melihat bagaimana identitas pengguna dapat menggeser proses penalaran, struktur respons, dan kualitas informasi yang diberikan model.

II.6 Penelitian Terdahulu dan Kesenjangan Penelitian

Pembahasan mengenai persona dan perilaku model bahasa telah menjadi bagian dari diskusi yang semakin luas dalam penelitian model berbasis transformator. Sejumlah studi sebelumnya memberikan gambaran mengenai bagaimana identitas pengguna, baik yang dinyatakan secara eksplisit maupun tersirat melalui gaya penulisan, dapat

memengaruhi respons model. Meskipun demikian, penelitian-penelitian tersebut umumnya memiliki cakupan yang terbatas pada satu jenis persona, satu kategori tugas, atau satu model tertentu. Bagian ini merangkum temuan utama dari penelitian terdahulu serta mengidentifikasi sejumlah kesenjangan yang melatarbelakangi penyusunan penelitian ini.

II.6.1 Ringkasan Literatur Terkait

Gupta et al. (Gupta dkk. 2024) menunjukkan bahwa persona eksplisit yang diberikan dalam instruksi dapat mengubah struktur penalaran model, bahkan ketika tugas yang diberikan tidak berkaitan dengan identitas sosial tersebut. Temuan ini membuka diskusi bahwa model tidak hanya memproses isi instruksi, tetapi juga memaknai identitas sebagai konteks tambahan yang membentuk langkah-langkah penalaran.

Di sisi lain, Tseng et al. (Tseng dkk. 2024) menyoroti fenomena persona implisit yang muncul dari gaya penulisan pengguna. Model dapat menafsirkan pilihan kata, tingkat formalitas, atau cara menyampaikan pertanyaan sebagai sinyal identitas, sehingga menghasilkan respons yang selaras dengan kategori sosial yang diasosiasikan dengan isyarat tersebut. Studi ini menunjukkan bahwa persona tidak harus dinyatakan secara eksplisit untuk memengaruhi respons model.

Penelitian lain menyoroti aspek ketidakstabilan penalaran model. Turpin et al. (Turpin dkk. 2023) menunjukkan bahwa perubahan kecil pada struktur instruksi dapat mengubah langkah *chain-of-thought* yang dihasilkan model. Hal ini menunjukkan bahwa proses penalaran model sangat bergantung pada konteks linguistik, termasuk gaya atau nada instruksi yang pada akhirnya berhubungan erat dengan persona.

Dalam konteks bias, Weidinger et al. (Weidinger dkk. 2021) menunjukkan bahwa model bahasa dapat memperkuat atau meniru pola stereotip yang ada dalam data pelatihan. Temuan ini relevan ketika dikaitkan dengan persona karena identitas pengguna dapat memperkuat pola bias tertentu, terutama pada tugas sosial dan moral yang melibatkan interpretasi nilai atau pengambilan posisi tertentu.

Penelitian mengenai personalisasi model, seperti yang dibahas dalam Naous et al. (Naous, Roziere, dkk. 2025), lebih menekankan bagaimana variasi preferensi pengguna dapat memengaruhi gaya atau struktur jawaban. Meskipun fokusnya berbeda, studi ini memberikan gambaran bahwa model bahasa memberikan respons yang bervariasi bergantung pada konteks identitas atau preferensi pengguna.

II.6.2 Keterbatasan Penelitian Sebelumnya

Meskipun penelitian-penelitian tersebut memberikan kontribusi penting dalam memahami hubungan antara persona dan perilaku model bahasa, sebagian besar studi masih memiliki sejumlah keterbatasan. Pertama, banyak penelitian hanya menggunakan satu model sehingga temuan yang diperoleh belum menggambarkan variasi perilaku antarmodel. Padahal, model yang berbeda dapat menunjukkan tingkat sensitivitas yang berbeda terhadap persona.

Kedua, cakupan persona yang diteliti cenderung terbatas, sering kali hanya mencakup beberapa persona eksplisit atau sejumlah contoh persona implisit yang relatif kecil. Kondisi ini membuat temuan penelitian sebelumnya belum cukup untuk menggambarkan bagaimana variasi persona yang lebih luas memengaruhi perilaku model.

Ketiga, sebagian besar penelitian hanya menguji satu atau dua jenis tugas. Padahal, persona dapat memengaruhi model secara berbeda pada penalaran numerik, penalaran logis, skenario sosial, maupun pertanyaan moral. Keterbatasan cakupan tugas ini membuat analisis sebelumnya belum mencerminkan penuh kompleksitas pengaruh persona.

Keempat, sebagian penelitian belum menyediakan kerangka evaluasi yang terotomatisasi dan konsisten. Tanpa mekanisme evaluasi yang terstruktur, sulit untuk memastikan bahwa perubahan respons benar-benar disebabkan oleh persona dan bukan oleh variasi lain seperti perbedaan prompt atau kejadian *drift* antarpernyataan.

II.6.3 Posisi dan Kontribusi Penelitian Ini

Penelitian ini disusun dengan mempertimbangkan keterbatasan-keterbatasan tersebut. Berbeda dengan penelitian sebelumnya, penelitian ini menggunakan pendekatan *multi model* dan *multi persona* untuk melihat bagaimana variasi identitas pengguna memengaruhi penalaran, gaya respons, dan kecenderungan bias. Dengan mengombinasikan beberapa kategori tugas—mulai dari penalaran numerik hingga skenario moral—penelitian ini bertujuan memberikan gambaran yang lebih utuh mengenai perilaku model bahasa ketika berinteraksi dengan berbagai persona pengguna.

Penelitian ini juga memanfaatkan *evaluation pipeline* yang terotomatisasi untuk memastikan konsistensi struktur instruksi dan mengurangi pengaruh variasi yang tidak diinginkan. Pendekatan ini diharapkan dapat memberikan hasil yang lebih stabil dan dapat direplikasi, sehingga memperkuat kontribusi penelitian dalam memahami

sensitivitas model bahasa terhadap persona.

Secara keseluruhan, penelitian-penelitian tersebut menunjukkan perlunya kajian yang lebih luas dan terstruktur mengenai bagaimana persona memengaruhi perilaku model bahasa, terutama ketika melibatkan lebih dari satu model dan lebih dari satu kategori tugas.

II.7 Penulisan Gambar, Tabel, Rumus, dan Kode

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

II.7.1 Gambar

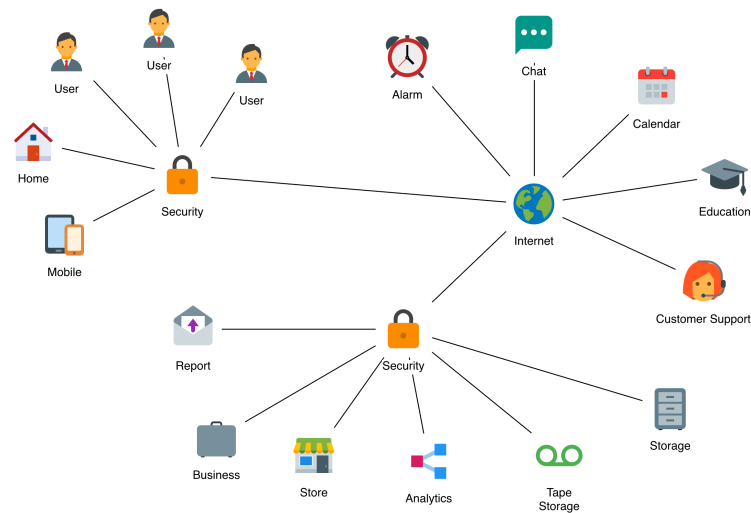
Contoh gambar dapat dilihat pada Gambar II.1. Gambar dan judulnya diposisikan di tengah. Nomor gambar tidak diakhiri tanda titik. Gambar tersebut dibuat menggunakan aplikasi draw.io dan disimpan ke format PNG setelah dengan zoom setting pada angka 300%. Ukuran gambar yang ditampilkan dapat diatur dengan mengubah nilai *width* dalam sintaks *includegraphics*.

Gambar umumnya tidak jelas atau kabur jika gambar tersebut:

- a. diperoleh dari hasil cropping pada suatu halaman buku atau situs web;
- b. hasil pembesaran gambar yang gambar aslinya sebenarnya berukuran kecil; atau
- c. disimpan dalam resolusi kecil

Ketidakjelasan gambar ini dapat dilihat pada garis-garis diagram yang tidak tegas dan tulisan-tulisan dalam gambar yang tampak kabur dan kurang jelas terbaca.

Untuk mendapatkan gambar yang tidak kabur (*blur*), langkah-langkah berikut dapat digunakan:



Gambar II.1 Contoh gambar jaringan

Tabel II.1 Tabel harga bahan pokok

Nama	Satuan	Harga
Buku	Exemplar	25000
Komputer	Unit	2500000
Pensil	Buah	118900

- Gambar yang didapat di suatu pustaka atau referensi sebaiknya digambar ulang, misalnya menggunakan PowerPoint, Canva, Figma, draw.io, atau yang lainnya.
- Jika diagram atau ilustrasi digambar menggunakan draw.io, saat gambar disimpan ke format PNG atau JPG (*export as*), lakukan *zoom* ke minimal 300% (*the default value is 100%*).
- Jika diagram digambar dengan menggunakan PowerPoint, gambar dapat langsung di-*copy-paste* ke Word.

II.7.2 Tabel

Tabel ada dua jenis, yaitu tabel yang bisa termuat dalam satu halaman dan tabel yang sangat panjang sehingga tidak muat dalam satu halaman.

II.7.2.1 Tabel yang Muat dalam Satu Halaman

Contoh tabel dapat dilihat pada Tabel II.1 dan II.2. Tabel dan judulnya dibuat rata kiri dan judul tabel diletakkan di atas tabel. Usahakan tabel dapat ditulis dalam satu halaman, tidak terpotong ke halaman berikutnya.

Tabel II.2 Tabel harga bahan sekunder

Nama	Satuan	Harga
Buku	Exemplar	25000
Komputer	Unit	2500000
Pensil	Buah	118900

II.7.2.2 Mengimpor Tabel dari Berkas Eksternal

Tabel II.3 diimpor dari berkas eksternal *table/tabel1.tex* menggunakan perintah *input*. Dengan demikian, jika tabel tersebut perlu diubah, cukup mengubah pada berkas eksternal tersebut tanpa perlu mengubah pada berkas utama ini.

Tabel II.3 Tabel harga bahan tertier

Nama	Satuan	Harga
Buku	Exemplar	25000
Komputer	Unit	2500000
Pensil	Buah	118900
Pensil	Buah	118900
Pensil	Buah	118900
Pensil	Buah	118900
Pensil	Buah	118900
Pensil	Buah	118900

II.7.2.3 Tabel yang Sangat Panjang

Jika tabel terlalu panjang sehingga tidak muat dalam satu halaman, gunakan paket *longtable* untuk membuat tabel yang dapat terpotong ke halaman berikutnya, seperti pada Tabel II.4.

Tabel II.4 Comprehensive Data Table Example

ID	Name	Score	Rank
1	Alice Smith	89	5
2	Bob Johnson	93	3
3	Carol Davis	95	2
4	Daniel Wilson	88	6
5	Eve Thompson	97	1
6	Frank Brown	85	7
7	Grace Lee	91	4
8	Henry Miller	80	9

Bersambung ke halaman berikutnya

Tabel II.4 Comprehensive Data Table Example (lanjutan)

ID	Name	Score	Rank
9	Irene Garcia	83	8
10	Jack Robinson	78	10
11	Kevin Harris	76	11
12	Laura Martin	75	12
13	Michael Clark	74	13
14	Natalie Lewis	73	14
15	Olivia Walker	72	15
16	Peter Hall	71	16
17	Quinn Allen	70	17
18	Rachel Young	69	18
19	Samuel King	68	19
20	Tina Wright	67	20
21	Uma Scott	66	21
22	Victor Green	65	22
23	Wendy Adams	64	23
24	Xavier Nelson	63	24
25	Yolanda Carter	62	25
26	Zachary Perez	61	26
27	Amelia Baker	60	27
28	Benjamin Rivera	59	28
29	Charlotte Rogers	58	29
30	David Murphy	57	30
31	Ethan Cooper	56	31
32	Fiona Reed	55	32
33	George Bailey	54	33
34	Hannah Cox	53	34
35	Isaac Howard	52	35
36	Julia Ward	51	36
37	Kyle Flores	50	37
38	Lily Bell	49	38
39	Mason Sanders	48	39
40	Nora Patterson	47	40
41	Owen Ramirez	46	41

Bersambung ke halaman berikutnya

Tabel II.4 Comprehensive Data Table Example (lanjutan)

ID	Name	Score	Rank
42	Penelope Torres	45	42
43	Quentin Foster	44	43
44	Rebecca Gonzales	43	44
45	Sebastian Bryant	42	45
46	Taylor Alexander	41	46
47	Ursula Russell	40	47
48	Vincent Griffin	39	48
49	William Diaz	38	49
50	Zoe Simmons	37	50

II.7.2.4 Beberapa Contoh Penulisan Rumus atau Persamaan Matematika Menggunakan LaTeX Termasuk Penomorannya

Contoh rumus matematika dapat ditulis seperti pada Persamaan II.1 di bawah ini. Penomoran persamaan diletakkan di sebelah kanan, dan rumus ditulis dalam mode *display math*.

$$E = mc^2 \quad (\text{II.1})$$

Contoh lain penulisan rumus matematika yang lebih kompleks dapat ditulis seperti pada Persamaan II.3.

$$f(x) = ax^2 + bx + c \quad (\text{II.2})$$

$$\begin{aligned} f'(x) &= \frac{d}{dx}(ax^2 + bx + c) \\ &= 2ax + b \end{aligned} \quad (\text{II.3})$$

Jika rumus terlalu panjang untuk ditulis dalam satu baris, gunakan lingkungan *multiline* seperti pada Persamaan II.4 di bawah ini.

$$\begin{aligned} y = a_0 + a_1x + a_2x^2 + a_3x^3 + a_4x^4 + a_5x^5 + a_6x^6 + a_7x^7 \\ + a_8x^8 + a_9x^9 + a_{10}x^{10} \end{aligned} \quad (\text{II.4})$$

Jika ada penurunan rumus yang terdiri dari beberapa baris, namun tidak memerlukan penomoran pada setiap baris, gunakan lingkungan *align**, misalnya:

$$\begin{aligned} S &= \sum_{i=1}^n i^2 \\ &= 1^2 + 2^2 + 3^2 + \cdots + n^2 \\ &= \frac{n(n+1)(2n+1)}{6} \end{aligned}$$

Contoh lainnya adalah rumus untuk mencari nilai rata-rata fungsi $f(x)$ pada interval $[p, q]$:

$$\begin{aligned} \bar{f} &= \frac{1}{q-p} \int_p^q f(x) dx \\ &= \frac{1}{q-p} \int_p^q (ax^2 + bx + c) dx \\ &= \frac{1}{q-p} \left[\frac{a}{3}x^3 + \frac{b}{2}x^2 + cx \right]_p^q \\ &= \frac{a(q^3 - p^3)}{3(q-p)} + \frac{b(q^2 - p^2)}{2(q-p)} + c \end{aligned}$$

II.7.3 Algoritma, Pseudocode, atau Kode

Contoh penulisan algoritma atau pseudocode dapat ditulis seperti pada Kode II.1 di bawah ini. Gunakan paket *listings* untuk menulis source code dalam bahasa pemrograman tertentu, seperti pada Kode II.2.

Kode II.1 Contoh pseudocode

```
ALGORITHM HelloWorld
    PRINT "Hello, World!"
END ALGORITHM
```

Kode II.2 Contoh source code Python

```
def hello_world():
    print("Hello, World!")
hello_world()
```

Tabel II.5 Contoh penggunaan kata "sedangkan" dan "sehingga"

Kata	Salah	Benar
sedangkan	Sedangkan sistem lama masih digunakan oleh banyak pengguna.	Sistem lama masih digunakan oleh banyak pengguna, sedangkan sistem baru belum siap.
sehingga	Sehingga sistem lama masih digunakan oleh banyak pengguna.	Sistem lama masih digunakan oleh banyak pengguna sehingga sistem baru belum siap.

II.8 Beberapa Kesalahan Penulisan yang Sering Terjadi

II.8.1 Penggunaan Kata "di mana" atau "dimana"

Banyak yang menuliskan kata "di mana" atau "dimana" sebagai pengganti kata "which" dalam bahasa Inggris. Padahal, penggunaan kata "di mana" atau "dimana" tidak tepat dalam konteks tersebut. Demikian juga untuk kata serupa, misalnya "yang mana". Kata "di mana" atau "dimana" ini harus diganti dengan kata lain, seperti "dengan", "tempat", "yang", dan sebagainya tergantung kalimatnya. Penjelasan lengkap dapat dilihat pada (BPBI).

II.8.2 Penggunaan Kata "sedangkan" dan "sehingga"

Kata "sedangkan" dan "sehingga" adalah kata hubung atau konjungsi. Konjungsi adalah kata atau ungkapan yang menghubungkan satuan bahasa (kata, frasa, klausa, dan kalimat). Konjungsi dapat dibagi menjadi konjungsi intrakalimat dan antarkalimat. Kata "sedangkan" menghubungkan dua klausa yang bersifat kontrasif, sedangkan "sehingga" menghubungkan dua klausa yang bersifat kausal. Dalam ragam formal, kata hubung "sedangkan" dan "sehingga" hanya dapat digunakan sebagai konjungsi intrakalimat sehingga kedua konjungsi itu **tidak dapat diletakkan pada awal kalimat**. Selain itu, penggunaan kata "sedangkan" harus didahului oleh koma (,), sedangkan kata "sehingga" tidak perlu didahului oleh koma (,). Contoh penggunaan yang benar dan salah dapat dilihat pada Tabel II.5.

II.8.3 Penggunaan Istilah yang Tidak Baku

Ada beberapa istilah yang sering digunakan dalam pembicaraan sehari-hari, tetapi tidak baku dalam penulisan ilmiah. Beberapa istilah tersebut antara lain:

1. analisa → analisis
2. eksisting atau existing → yang ada atau saat ini
3. bisnis proses → proses bisnis

4. user → pengguna
5. system → sistem
6. database → basis data
7. aktifitas → aktivitas
8. efektifitas → efektivitas
9. sosial media → media sosial

II.8.4 Pemisah Desimal dan Ribuan

Tanda pemisah desimal dalam bahasa Indonesia adalah tanda koma, contoh:

1. (Salah) Akurasi naik menjadi 50.6%
2. (Benar) Akurasi naik menjadi 50,6%

II.8.5 Daftar atau *List*

Ada beberapa aturan penulisan daftar atau *list* yang perlu diperhatikan, antara lain:

- a) Jika memungkinkan, hindari penggunaan “bullet points” atau sejenisnya. Sebaiknya, gunakan angka (1, 2, 3, ...) atau huruf (a, b, c, ...). Dengan demikian, pembaca dapat dengan mudah melihat jumlah *item* atau *list*.
- b) Jika dalam daftar hanya ada satu item, tidak perlu menggunakan nomor urut.
- c) Penjelasan atau deskripsi suatu item sebaiknya menyatu dengan judul item tersebut, tidak berbeda halaman. Contoh yang salah: judul item ada di halaman 10, namun deskripsinya di halaman 11. Sebaiknya pindahkan judul tersebut ke halaman 11.
- d) Jika penjelasan atau deskripsi suatu item cukup panjang, misalnya lebih dari 1 halaman atau terdiri atas beberapa paragraf, sebaiknya setiap item tersebut dijadikan judul subbab, kecuali jika level subbab sudah mencapai level 4.

II.8.6 Penggunaan Kata “masing-masing” dan “setiap”

Kata “masing-masing” digunakan di belakang kata yang diterangkan, misalnya “Setiap proses menggunakan algoritma masing-masing”. Kata “tiap-tiap” atau “setiap” ditempatkan di depan kata yang diterangkan, misalnya “Setiap proses menggunakan algoritma tertentu”.

BAB III

ANALISIS MASALAH

III.1 Analisis Kondisi Saat Ini

Menurut **laudon2020**<empty citation>, gambarkan terlebih dahulu model konseptual sistem yang ada saat ini. Model konseptual ini berisi berbagai komponen atau subsistem dan interaksi antarsubsistem tersebut. Setelah itu, berikan penjelasan tentang masalah yang ada pada sistem tersebut. Paragraf berikut berisi contoh penjabaran masalah sistem informasi fasilitas kesehatan untuk pasien (**pressman2019**).

III.2 Analisis Kebutuhan

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

III.2.1 Identifikasi Masalah Pengguna

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetur.

III.2.2 Kebutuhan Fungsional

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

III.2.3 Kebutuhan Nonfungsional

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

III.3 Analisis Pemilihan Solusi

III.3.1 Alternatif Solusi

Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Donec odio elit, dictum in, hendrerit sit amet, egestas sed, leo. Praesent feugiat sapien aliquet odio. Integer vitae justo. Aliquam vestibulum fringilla lorem. Sed neque lectus, consectetur at, consectetur sed, eleifend ac, lectus. Nulla facilisi. Pellentesque eget lectus. Proin eu metus. Sed porttitor. In hac habitasse platea dictumst. Suspendisse eu lectus. Ut mi mi, lacinia sit amet, placerat et, mollis vitae, dui. Sed ante tellus, tristique ut, iaculis eu, malesuada ac, dui. Mauris nibh leo, facilisis non, adipiscing quis, ultrices a, dui.

III.3.2 Analisis Penentuan Solusi

Morbi luctus, wisi viverra faucibus pretium, nibh est placerat odio, nec commodo wisi enim eget quam. Quisque libero justo, consectetur a, feugiat vitae, porttitor eu, libero. Suspendisse sed mauris vitae elit sollicitudin malesuada. Maecenas ultricies eros sit amet ante. Ut venenatis velit. Maecenas sed mi eget dui varius euismod.

Phasellus aliquet volutpat odio. Vestibulum ante ipsum primis in faucibus orci luctus et ultrices posuere cubilia Curae; Pellentesque sit amet pede ac sem eleifend consetetuer. Nullam elementum, urna vel imperdiet sodales, elit ipsum pharetra ligula, ac pretium ante justo a nulla. Curabitur tristique arcu eu metus. Vestibulum lectus. Proin mauris. Proin eu nunc eu urna hendrerit faucibus. Aliquam auctor, pede consequat laoreet varius, eros tellus scelerisque quam, pellentesque hendrerit ipsum dolor sed augue. Nulla nec lacus.

BAB IV

DESAIN KONSEP SOLUSI

Ilustrasikan desain konsep solusi dalam bentuk model konseptual dan penjelasan secara ringkas, beserta perbedaannya dengan sistem saat ini. Ilustrasi harus dapat dibandingkan (*before and after*). Karena masih berupa proposal, bab ini hanya berisi gambar desain konsep solusi tersebut dan penjelasan perbandingannya dengan gambar sistem yang ada saat ini (yang tergambar di awal Bab III).

BAB V

RENCANA SELANJUTNYA

Jelaskan secara detail langkah-langkah rencana selanjutnya, hal-hal yang diperlukan atau akan disiapkan, dan risiko dan mitigasinya, yang meliputi:

1. Rencana implementasi, termasuk alat dan bahan yang diperlukan, lingkungan, konfigurasi, biaya, dan sebagainya.
2. Desain pengujian dan evaluasi, misalnya metode verifikasi dan validasi.
3. Analisis risiko dan mitigasi, misalnya tindakan selanjutnya jika ada yang tidak berjalan sesuai rencana.

DAFTAR PUSTAKA

- Bommasani, Rishi, Drew A. Hudson, Ehsan Adeli, dkk. 2021. “On the Opportunities and Risks of Foundation Models”. *arXiv preprint arXiv:2108.07258*, <https://arxiv.org/abs/2108.07258>.
- Gupta, Shashank, Vaishnavi Shrivastava, Ameet Deshpande, Ashwin Kalyan, Peter Clark, Ashish Sabharwal, dan Tushar Khot. 2024. “Bias Runs Deep: Implicit Reasoning Biases in Persona-Assigned Language Models”. Dalam *Proceedings of the Twelfth International Conference on Learning Representations*. <https://openreview.net/forum?id=kGteeZ18Ir>.
- Naous, Tarek, Baptiste Roziere, dkk. 2025. “Training and Evaluating User Language Models”. *arXiv preprint arXiv:2510.06552*, <https://arxiv.org/abs/2510.06552>.
- Tseng, Yu-Min, Yu-Chao Huang, Teng-Yun Hsiao, Wei-Lin Chen, Chao-Wei Huang, Yu Meng, dan Yun-Nung Chen. 2024. “Two Tales of Persona in LLMs: A Survey of Role-Playing and Personalization”. Dalam *Findings of the Association for Computational Linguistics: EMNLP 2024*, 16612–16631. <https://aclanthology.org/2024.findings-emnlp.969>.
- Turpin, Miles, dkk. 2023. “Language Models Don’t Always Say What They Think: Unfaithful Explanations in Chain-of-Thought Reasoning”. *arXiv preprint arXiv:2305.04388*, <https://arxiv.org/abs/2305.04388>.
- Weidinger, Laura, John Mellor, Maribeth Rauh, Christopher Griffin, Iason Gabriel, Jonathan Uesato, Po-Sen Huang, Zachary Kenton, Tom B. Brown, dkk. 2021. “Ethical and Social Risks of Harm from Language Models”. *arXiv preprint arXiv:2112.04359*, <https://arxiv.org/abs/2112.04359>.
- Zhao, Yanhao, Eric Wallace, Shi Feng, Mohit Singh, dan Matt Gardner. 2021. “Calibrate Before Use: Improving Few-Shot Performance of Language Models”. Dalam *Proceedings of the International Conference on Machine Learning*, 12697–12706.

Zhou, Luozhi, dkk. 2023. “Large Language Models Are Sensitive to Prompt Framing”. *arXiv preprint arXiv:2310.05400*, <https://arxiv.org/abs/2310.05400>.