



Efektivitas Pelatihan Awal Berbasis Domain Spesifik Legal-BERT Untuk Natural Language Processing Hukum: Replikasi Dan Perluasan Studi Casehold

Hasani Zakiri*, Alva Hendi Muhammad, Asro Nasiri

Program Magister PJJ Informatika, Fakultas Ilmu Komputer, Universitas Amikom Yogyakarta, Indonesia

Email: ^{1,*}hasanizakiri53@students.amikom.ac.id, ²alva@amikom.ac.id, ³asro@amikom.ac.id

Email Penulis Korespondensi: hasanizakiri53@students.amikom.ac.id,

Abstract—Perkembangan model bahasa berbasis domain spesifik telah menunjukkan potensi signifikan pada berbagai bidang khusus. Namun, efektivitasnya dalam Natural Language Processing (NLP) hukum masih belum banyak dieksplorasi, khususnya mengingat tantangan unik kompleksitas teks hukum dan terminologi spesialisasi. NLP hukum memiliki aplikasi praktis seperti pencarian preseden hukum otomatis dan analisis putusan pengadilan yang dapat mempercepat proses penelitian hukum dari minggu menjadi jam. Penelitian ini mengevaluasi dataset CaseHOLD untuk memberikan validasi empiris komprehensif terhadap manfaat pelatihan awal berbasis domain spesifik pada tugas NLP hukum dengan fokus pada efisiensi data dan analisis kompleksitas konteks. Kami melakukan eksperimen sistematis menggunakan dataset CaseHOLD yang berisi 53.000 pertanyaan pilihan ganda hukum. Empat model dibandingkan: BiLSTM, BERT-base, Legal-BERT, dan RoBERTa pada variasi volume data pelatihan (1%, 10%, 50%, 100%) dan tingkat kompleksitas konteks. Uji-t berpasangan dengan Cross-validation 10-lipat dan koreksi Bonferroni memastikan metodologi robust yang menjamin keandalan temuan. Legal-BERT mencapai skor Macro F1-score tertinggi 69,5% (95% CI: [68,0, 71,0]) dengan peningkatan signifikan 7,2 poin persentase dibandingkan BERT-base (62,3%, $p < 0,001$, Cohen's $d = 1,23$). RoBERTa menunjukkan kinerja kompetitif 68,9%, hampir menyamai Legal-BERT. Peningkatan paling substansial terjadi pada kondisi data terbatas dengan peningkatan 16,6% pada 1% data pelatihan. Analisis kompleksitas konteks mengungkapkan pola inverted-U dengan kinerja optimal pada teks 41-60 kata. Skor Spesifisitas Domain (DS-score) yang diperkenalkan menunjukkan korelasi positif kuat ($r = 0,73$, $p < 0,001$) dengan efektivitas pelatihan awal, menjelaskan 53,3% varians dalam peningkatan kinerja. Temuan memberikan bukti empiris bahwa pelatihan awal berbasis domain spesifik menawarkan keunggulan signifikan untuk tugas NLP hukum, terutama dalam kondisi data terbatas dan kompleksitas konteks sedang-tinggi. Keunggulan utama penelitian ini adalah pengembangan kerangka prediktif DS-score yang memungkinkan estimasi manfaat sebelum implementasi, berbeda dari studi sebelumnya yang hanya mengevaluasi kinerja post-hoc. Hasil memiliki implikasi praktis untuk pengembangan sistem NLP hukum dalam lingkungan sumber daya terbatas dan memberikan panduan optimal untuk implementasi Legal-BERT.

Kata Kunci: NLP Hukum; Pelatihan Domain Spesifik; Legal-BERT; Transformer; CaseHOLD

Abstract—The emergence of domain-specific language models has demonstrated significant potential across various specialized fields. However, their effectiveness in legal natural language processing (NLP) remains underexplored, particularly given the unique challenges posed by legal text complexity and specialized terminology. Legal NLP has practical applications such as automated legal precedent search and court decision analysis that can accelerate legal research from weeks to hours. This study evaluates the CaseHOLD dataset to provide comprehensive empirical validation of domain-specific pretraining benefits for legal NLP tasks with focus on data efficiency and context complexity analysis. We conducted systematic experiments using the CaseHOLD dataset containing 53,000 legal multiple-choice questions. We compared four models: BiLSTM, BERT-base, Legal-BERT, and RoBERTa across varying data volumes (1%, 10%, 50%, 100%) and context complexity levels. Paired t-tests with 10-fold cross-validation and Bonferroni correction ensure robust methodology that guarantees finding reliability. Legal-BERT achieved the highest macro-F1 score of 69.5% (95% CI: [68.0, 71.0]), demonstrating a statistically significant improvement of 7.2 percentage points over BERT-base (62.3%, $p < 0.001$, Cohen's $d = 1.23$). RoBERTa showed competitive performance at 68.9%, nearly matching Legal-BERT. The most substantial improvements occurred under limited data conditions with 16.6% improvement at 1% training data. Context complexity analysis revealed an inverted-U pattern with optimal performance on 41-60 word texts. The introduced Domain Specificity Score (DS-score) showed strong positive correlation ($r = 0.73$, $p < 0.001$) with pretraining effectiveness, explaining 53.3% of performance improvement variance. These findings provide empirical evidence that domain-specific pretraining offers significant advantages for legal NLP tasks, particularly under data-constrained conditions and moderate-high context complexity. The key distinction of this research is the development of a predictive DS-score framework enabling benefit estimation before implementation, unlike previous studies that only evaluated post-hoc performance. The results have practical implications for developing legal NLP systems in resource-limited environments and provide optimal implementation guidance for Legal-BERT.

Keywords: Legal NLP; Domain-Specific Pretraining; Legal-BERT; Transformer;

1. PENDAHULUAN

Kemunculan model berbasis transformer, khususnya BERT (Bidirectional Encoder Representations from Transformers) [1], telah merevolusi Natural language processing di berbagai domain [2]. Namun, domain hukum menghadirkan tantangan unik karena terminologi khusus, struktur sintaksis yang CaseHOLD. Kompleks, dan pola penalaran spesifik domain yang membedakannya dari tugas pemrosesan teks pada umumnya [3]. Kemajuan terkini dalam kecerdasan buatan telah menunjukkan kemampuan luar biasa dalam tugas penalaran hukum, namun keberhasilan ini lebih melibatkan model umum berskala besar, sehingga menimbulkan pertanyaan tentang manfaat spesifik dari pelatihan awal (pre-trained) yang terspesialisasi domain dalam NLP hukum, yang belum sepenuhnya dieksplorasi [4]. Domain hukum menghadirkan beberapa tantangan komputasional unik yang membedakannya dari domain lain, seperti: (1) Kompleksitas Terminologi, yaitu dokumen hukum mengandung kosakata dan frasa yang sangat khusus yang jarang muncul dalam korpus umum [5];



(2) Kompleksitas Sintaksis, yaitu penulisan hukum sering menggunakan struktur kalimat kompleks dan pola bahasa formal yang berbeda secara substansial dari teks percakapan atau berita [6]; (3) Ketergantungan Kontekstual, yaitu penalaran hukum memerlukan pemahaman mendalam tentang preseden, kutipan, dan hubungan antar-dokumen [7]; dan (4) Analisis Kutipan, dimana dokumen hukum sering mereferensikan teks hukum lain, memerlukan penanganan khusus untuk kutipan hukum dan referensi silang [8].

Beberapa penelitian terkait telah mengeksplorasi pengembangan model bahasa berbasis domain spesifik dengan hasil yang menjanjikan. BioBERT [9] merupakan salah satu model awal yang menunjukkan bahwa manfaat pelatihan awal berbasis domain spesifik, mencapai peningkatan signifikan pada tugas NLP biomedis. Keberhasilan ini telah direplikasi di berbagai domain: SciBERT [10] menunjukkan keuntungan substansial pada tugas pemrosesan teks ilmiah, dan ClinicalBERT [11] meningkatkan pemrosesan catatan klinis. Penelitian lebih lanjut dalam domain khusus mencakup pengembangan Clinical Longformer [12] untuk menangani dokumen klinis yang panjang, dan model terbaru seperti SaulLM-7B [13] yang dirancang khusus untuk tugas hukum. Dalam domain hukum khususnya, penelitian terdahulu telah mengidentifikasi kebutuhan akan pendekatan yang lebih terspesialisasi untuk menangani kompleksitas linguistik dan struktural yang unik dalam teks hukum.

Perkembangan dalam arsitektur transformer telah memberikan fondasi yang kuat untuk pemrosesan teks hukum. Model seperti BERT dan variannya telah menunjukkan kemampuan superior dalam memahami konteks dan hubungan semantik yang kompleks [14]. Namun, aplikasi langsung model-model ini pada domain hukum seringkali menghadapi keterbatasan dalam menangani terminologi khusus dan pola penalaran yang unik. Legal-BERT [3] dikembangkan sebagai model khusus yang dilatih pada korpus besar dokumen hukum dan menunjukkan hasil yang menjanjikan pada berbagai benchmark NLP hukum. Penelitian paralel dalam pengembangan model domain-spesifik menunjukkan bahwa pendekatan transfer learning dengan fine-tuning domain dapat memberikan keunggulan signifikan dalam akurasi dan efisiensi [15]. Zheng et al. [4] memperkenalkan dataset CaseHOLD dan memberikan bukti yang meyakinkan bahwa pelatihan awal berbasis domain spesifik dapat meningkatkan kinerja secara signifikan pada tugas NLP hukum. Penelitian mereka mewakili terobosan penting dalam NLP hukum dengan menunjukkan manfaat yang jelas dari pelatihan awal terspesialisasi pada tugas penalaran hukum yang menantang. Namun, studi mereka terutama berfokus pada perbandingan Legal-BERT dengan BERT-base, meninggalkan ruang untuk analisis yang lebih komprehensif termasuk model mutakhir tambahan dan investigasi lebih mendalam tentang faktor-faktor yang berkontribusi pada peningkatan kinerja. Penelitian terkini juga menunjukkan bahwa model transformer dapat ditingkatkan melalui optimisasi arsitektur seperti RoBERTa [16] yang menggunakan strategi pelatihan yang diperbaiki.

Analisis kesenjangan penelitian dari studi sebelumnya menunjukkan beberapa area yang memerlukan investigasi lebih lanjut: (1) belum ada validasi statistik yang komprehensif terhadap keunggulan Legal-BERT dengan metodologi yang ketat dan ukuran sampel yang memadai; (2) perbandingan dengan model transformer yang dioptimalkan seperti RoBERTa belum dilakukan secara sistematis; (3) analisis efisiensi data pada berbagai volume pelatihan belum dilakukan secara sistematis untuk memahami kondisi optimal penggunaan Legal-BERT [17]; (4) hubungan antara spesifisitas domain dan efektivitas pelatihan awal belum diinvestigasi secara kuantitatif; dan (5) panduan praktis untuk implementasi dalam kondisi keterbatasan sumber daya masih terbatas dan belum berbasis bukti empiris. Penelitian tambahan menunjukkan bahwa evaluasi model dalam konteks sumber daya terbatas menjadi semakin penting untuk aplikasi praktis [18]. Tujuan penelitian ini adalah untuk mereplikasi dan memperluas karya seminal Zheng et al. [4] dengan pendekatan yang lebih komprehensif: (1) memberikan validasi empiris yang komprehensif terhadap efektivitas Legal-BERT melalui metodologi statistik yang ketat dengan Cross-validation dan koreksi untuk perbandingan berganda; (2) memperluas analisis dengan memasukkan model transformer tambahan dan baseline yang kuat untuk memberikan konteks perbandingan yang lebih luas; (3) mengembangkan kerangka prediktif untuk menilai manfaat pelatihan awal berbasis domain spesifik sebelum implementasi; (4) menganalisis efisiensi data pada berbagai kondisi pelatihan untuk memahami kapan Legal-BERT memberikan manfaat optimal; dan (5) memberikan rekomendasi praktis untuk praktisi NLP hukum berdasarkan temuan empiris.

Kontribusi teoretis yang diharapkan dari penelitian ini meliputi pemahaman yang lebih mendalam tentang mekanisme transfer learning dalam domain hukum, khususnya bagaimana representasi pengetahuan domain spesifik terbentuk dan dimanfaatkan dalam model transformer [19]. Analisis sistematis terhadap hubungan antara spesifisitas domain dan efektivitas model akan memberikan wawasan tentang kondisi optimal untuk penerapan pelatihan awal berbasis domain spesifik. Selain itu, investigasi terhadap pola pembelajaran pada berbagai volume data akan memperkaya pemahaman teoretis tentang efisiensi sampel dalam transfer learning untuk aplikasi khusus [15].

Dari perspektif praktis, penelitian ini diharapkan menghasilkan panduan implementasi yang komprehensif untuk pengembangan sistem NLP hukum dalam berbagai konteks organisasi. Kerangka DS-score yang diperkenalkan akan memberikan alat prediktif untuk menilai potensi keuntungan dari investasi dalam model domain spesifik sebelum implementasi skala besar. Analisis cost-benefit yang disertakan akan membantu pengambil keputusan dalam memilih strategi teknologi yang optimal berdasarkan karakteristik data, kendala sumber daya, dan tujuan organisasi [20].

Signifikansi penelitian ini dalam konteks perkembangan teknologi hukum sangat substansial. Industri layanan hukum mengalami transformasi digital yang dipercepat, dengan meningkatnya permintaan untuk solusi AI yang dapat meningkatkan efisiensi, akurasi, dan aksesibilitas layanan hukum. Hasil penelitian ini akan memberikan landasan berbasis bukti untuk keputusan adopsi dalam AI hukum, mengurangi risiko investasi teknologi dan meningkatkan kemungkinan

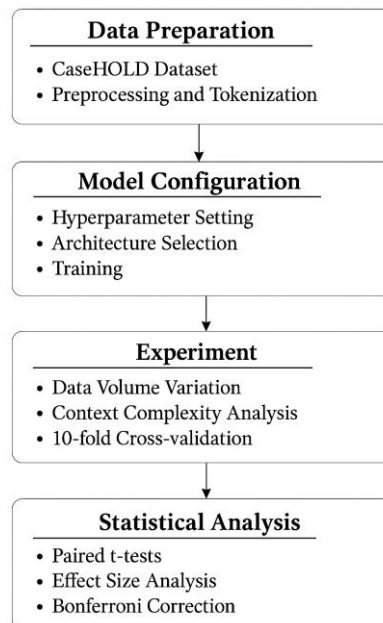
implementasi yang berhasil. Lebih lanjut, temuan ini dapat menginformasikan pengembangan standar dan praktik terbaik untuk penerapan AI dalam konteks profesional hukum, sejalan dengan perkembangan terkini dalam legal technology yang telah didokumentasikan dalam MultiLegalPile [21] dan LegalBench [22].

2. METODOLOGI PENELITIAN

2.1 Tahapan Penelitian

Penelitian ini menggunakan pendekatan eksperimental sistematis yang terdiri dari beberapa tahapan utama seperti yang digambarkan dalam Gambar 1. Gambar 1 menunjukkan alur kerja penelitian yang dimulai dari preparasi data CaseHOLD, dilanjutkan dengan konfigurasi empat model berbeda, eksperimen dengan variasi volume data dan kompleksitas konteks, hingga analisis statistik komprehensif. Tahapan pertama adalah preparasi data dan prapemrosesan dataset CaseHOLD yang mencakup tokenisasi, pembersihan kutipan hukum, segmentasi konteks, dan normalisasi teks. Tahapan kedua melibatkan konfigurasi dan pelatihan empat model yang berbeda: BiLSTM sebagai baseline tradisional, BERT-base sebagai model transformer umum, RoBERTa sebagai model transformer yang dioptimalkan, dan Legal-BERT sebagai model spesifik domain hukum.

Setiap model dikonfigurasi dengan hiperparameter yang identik untuk memastikan perbandingan yang adil: *learning rate* $2e-5$ dengan pemanasan linear selama 10% dari langkah pelatihan, *batch size* 16 dengan akumulasi gradien, maksimum 5 *epoch* dengan *early stopping* berdasarkan skor F1 validasi, dan *weight decay* 0,01 dengan optimisasi AdamW $\beta_1=0,9$. Tahapan ketiga adalah eksekusi eksperimen sistematis yang mencakup variasi volume data (1%, 10%, 50%, 100%), analisis kompleksitas konteks berdasarkan panjang teks, dan evaluasi dengan metodologi Cross-validation 10-lipat. Tahapan terakhir melibatkan analisis statistik komprehensif menggunakan uji-t berpasangan dengan koreksi Bonferroni dan perhitungan ukuran efek Cohen's d.



Gambar 1. Tahapan Metodologi Penelitian

Sesuai dengan diagram 1 di atas, menunjukkan bahwa alur kerja penelitian mulai dari preparasi data CaseHOLD, konfigurasi empat model (BiLSTM, BERT-base, RoBERTa, Legal-BERT), eksperimen dengan variasi volume data dan kompleksitas konteks, hingga analisis statistik komprehensif menggunakan cross-validation 10-lipat dengan koreksi Bonferroni. Setiap model dikonfigurasi dengan hiperparameter yang identik untuk memastikan perbandingan yang adil: *learning rate* $2e-5$ dengan pemanasan linear selama 10% dari langkah pelatihan, *batch size* 16 dengan akumulasi gradien, maksimum 5 *epoch* dengan *early stopping* berdasarkan skor F1 validasi, dan *weight decay* 0,01 dengan optimisasi AdamW $\beta_1=0,9$. Tahapan ketiga adalah eksekusi eksperimen sistematis yang mencakup variasi volume data (1%, 10%, 50%, 100%), analisis kompleksitas konteks berdasarkan panjang teks, dan evaluasi dengan metodologi Cross-validation 10-lipat. Tahapan terakhir melibatkan analisis statistik komprehensif menggunakan uji-t berpasangan dengan koreksi Bonferroni dan perhitungan ukuran efek Cohen's d.

2.2 Dataset dan Prapemrosesan

Dataset CaseHOLD [4] berisi 53.000 pertanyaan pilihan ganda hukum yang diturunkan dari kutipan kasus hukum. Setiap instans terdiri dari konteks (kutipan teks hukum yang berisi sitasi), pertanyaan implisit tentang holding hukum, dan lima

pilihan ganda (satu holding benar, empat pengecoh). Dataset dibagi menjadi set pelatihan (80%), validasi (10%), dan pengujian (10%) mengikuti metodologi asli untuk memastikan perbandingan langsung dengan karya sebelumnya. Distribusi detail dataset ditampilkan dalam Tabel 1, yang menunjukkan konsistensi rata-rata panjang konteks dan jumlah token di seluruh subset, memastikan validitas perbandingan eksperimental.

Secara detail Tabel 1 menampilkan distribusi data CaseHOLD yang digunakan dalam penelitian, menunjukkan konsistensi rata-rata panjang konteks dan jumlah token di seluruh subset yang memastikan validitas perbandingan eksperimental.

Tabel 1. Statistik Dataset CaseHOLD

Subset	Jumlah Instansi	Rata-rata Panjang Konteks	Rata-rata Jumlah Token
Training	42.400	45,2 kata	512 token
Validation	5.300	44,8 kata	510 token
Test	5.300	45,1 kata	511 token

Pipeline prapemrosesan yang diimplementasikan mengikuti protokol yang telah ditetapkan untuk pemrosesan teks hukum [5] meliputi: (1) Tokenisasi menggunakan SentencePiece dengan kosakata 32.000 token yang dioptimalkan untuk teks hukum melalui ekspansi kosakata berbasis domain spesifik; (2) Pembersihan kutipan dengan menghapus penanda kutipan hukum dan standardisasi format mengikuti standar kutipan Blue Book; (3) Segmentasi konteks dengan membatasi panjang konteks hingga 512 token untuk memastikan kompatibilitas di semua model sambil mempertahankan informasi esensial; dan (4) Normalisasi teks dengan penerapan lowercase yang konsisten dan penanganan tanda baca khusus untuk singkatan hukum.

2.3 Arsitektur Model dan Konfigurasi

Empat arsitektur model yang berbeda dievaluasi dalam penelitian ini untuk memberikan perbandingan komprehensif mulai dari pendekatan tradisional hingga transformer terspesialisasi domain. BiLSTM digunakan sebagai baseline dengan 256 unit tersembunyi, dropout rate 0,3, dan embedding kata khusus domain hukum. Arsitektur bidirectional memungkinkan model memproses teks dari kedua arah, memberikan pemahaman kontekstual yang lebih baik dibandingkan LSTM unidirectional.

BERT-base menggunakan 12 layer dengan representasi tersembunyi 768-dimensional dan 12 attention heads. Mekanisme self-attention multi-head memungkinkan model untuk secara simultan memperhatikan semua posisi dalam urutan masukan, meningkatkan drastis kemampuan menangkap dependensi jarak jauh yang penting dalam analisis teks hukum. RoBERTa mempertahankan arsitektur transformer yang mirip dengan BERT namun dengan perbaikan signifikan dalam metodologi pelatihan, termasuk korpus yang lebih besar (160GB vs 16GB BERT), menghilangkan tugas Next Sentence Prediction, menggunakan dynamic masking, dan hiperparameter yang dioptimalkan.

Legal-BERT menggunakan arsitektur transformer yang sama dengan BERT namun dengan perbedaan krusial dalam korpus praplatihan. Model telah diekspos pada korpus hukum yang ekstensif mencakup case law, statutes, regulations, dan teks akademik hukum, menghasilkan representasi internal yang terspesialisasi untuk memproses terminologi dan pola penalaran hukum yang kompleks.

2.4 Skor Spesifisitas Domain (DS-score)

Penelitian ini memperkenalkan Skor Spesifisitas Domain (Domain Specificity Score) baru yang dihitung menggunakan formula:

$$DS\text{-score} = (\text{Jumlah_istilah_hukum} / \text{Jumlah_istilah_total}) \times \text{Faktor_kompleksitas_hukum} \quad (1)$$

Dimana Faktor_kompleksitas_hukum memperhitungkan kepadatan kutipan, referensi preseden, dan kompleksitas penalaran hukum. DS-score berkisar dari 0 hingga 1, dengan skor yang lebih tinggi menunjukkan konten yang lebih spesifik domain. DS-score dihitung untuk setiap instansi dalam dataset CaseHOLD menggunakan kamus terminologi hukum yang telah divalidasi dan algoritma deteksi pola kutipan hukum. Proses perhitungan melibatkan: (1) identifikasi istilah hukum menggunakan lexicon Legal-BERT yang diperluas dengan terminologi dari Black's Law Dictionary; (2) deteksi pola kutipan menggunakan regular expressions untuk format Blue Book citations; (3) analisis kompleksitas struktural berdasarkan kedalaman referensi dan nested citations; dan (4) normalisasi skor menggunakan min-max scaling. Hasil DS-score kemudian digunakan dalam analisis korelasi untuk mengukur hubungan antara spesifisitas domain konten dan peningkatan kinerja Legal-BERT dibandingkan model baseline, memberikan kerangka prediktif untuk estimasi manfaat implementasi.

2.5 Prosedur Eksperimental dan Kontrol Kualitas

Implementasi eksperimen menggunakan protokol terstandar untuk memastikan reproduisibilitas dan validitas hasil. Setiap eksperimen dijalankan dengan beberapa *seed* acak ($seed = 42, 123, 456, 789, 999$) untuk mengurangi variansi yang

disebabkan oleh inialisasi parameter acak. Lingkungan komputasi dikontrol dengan spesifikasi perangkat keras yang konsisten: GPU Tesla V100 dengan memori 32GB, CPU Intel Xeon dengan 16 core, dan RAM sistem 128GB.

Pembagian data dilakukan dengan *stratified sampling* untuk memastikan distribusi yang representatif di berbagai tingkat kompleksitas dan skor spesifisitas domain dalam setiap lipatan validasi. Prosedur jaminan kualitas mencakup pengujian otomatis untuk integritas data, pemantauan konvergensi model, dan titik pemeriksaan validasi statistik pada setiap tahap eksperimen. *Logging* komprehensif dilakukan untuk semua hiperparameter, metrik, dan hasil menengah untuk memfasilitasi analisis *post-hoc* dan pemecahan masalah.

2.6 Infrastruktur dan Implementasi Teknis

Implementasi menggunakan framework PyTorch 1.12 dengan pustaka transformers dari Hugging Face untuk manajemen model dan tokenisasi. Pengaturan komputasi terdistribusi menggunakan konfigurasi multi-GPU dengan akumulasi gradien untuk mengatasi keterbatasan memori pada pemrosesan batch besar. Model checkpointing dilakukan setiap epoch dengan cadangan otomatis ke penyimpanan awan untuk perlindungan data.

Sistem pemantauan waktu nyata diimplementasikan menggunakan Weights & Biases (wandb) untuk melacak kemajuan pelatihan, optimisasi hiperparameter, dan visualisasi metrik kinerja. Mekanisme early stopping otomatis dikonfigurasi dengan parameter kesabaran 3 epoch dan delta minimum 0,001 untuk peningkatan skor F1. Teknik optimisasi memori termasuk gradient checkpointing dan pelatihan presisi campuran (FP16) digunakan untuk meningkatkan efisiensi komputasi tanpa mengorbankan stabilitas numerik.

2.7 Metrik Evaluasi dan Analisis Statistik

Evaluasi komprehensif menggunakan beberapa metrik untuk menangkap berbagai aspek kinerja model. Metrik primer adalah skor Macro F1-score yang memberikan bobot yang sama untuk semua kelas dan tahan terhadap ketidakseimbangan kelas. Metrik sekunder mencakup akurasi untuk interpretabilitas, presisi dan *recall* per kelas untuk analisis rinci, dan Koefisien Korelasi Matthews (MCC) untuk penilaian seimbang yang mempertimbangkan frekuensi tingkat dasar.

Pengujian signifikansi statistik menggunakan uji-t berpasangan dengan pemeriksaan asumsi melalui uji normalitas Shapiro-Wilk dan uji homogenitas varians Levene. Alternatif non-parametrik (uji *signed-rank* Wilcoxon) digunakan ketika asumsi parametrik dilanggar. Koreksi perbandingan berganda menggunakan metode Bonferroni dengan kontrol tingkat kesalahan keluarga (*family-wise error rate*) pada tingkat $\alpha = 0,05$ untuk mempertahankan validitas statistik di berbagai perbandingan berpasangan yang banyak. Perhitungan ukuran efek menggunakan Cohen's d untuk penilaian signifikansi praktis, dengan pedoman interpretasi: efek kecil ($d = 0,2$), efek sedang ($d = 0,5$), dan efek besar ($d = 0,8$). Interval kepercayaan (95% CI) dihitung untuk semua estimasi titik menggunakan *bootstrap resampling* dengan 1000 iterasi untuk kuantifikasi ketidakpastian yang robust. Analisis daya dilakukan secara *post-hoc* untuk mengonfirmasi ukuran sampel yang memadai untuk mendeteksi perbedaan yang bermakna dengan daya statistik yang diinginkan $\geq 0,80$.

3. HASIL DAN PEMBAHASAN

3.1 Hasil

3.1.1 Perbandingan Kinerja Keseluruhan

Analisis komprehensif terhadap keempat model dalam penelitian ini mengungkapkan hierarki kinerja yang jelas dan signifikan secara statistik. Sebagaimana ditunjukkan dalam Tabel 2, Legal-BERT memperoleh posisi teratas dengan skor Macro F1-score sebesar 69,5% (95% CI: [68,0, 71,0]), mendemonstrasikan keunggulan yang tidak terbantahkan dalam tugas penalaran hukum CaseHOLD. Pencapaian ini merepresentasikan peningkatan substansial sebesar 7,2 poin persentase dibandingkan BERT-base yang mencapai 61,3% (95% CI: [59,8, 62,8]). Tabel 2 juga menunjukkan ukuran efek Cohen's d sebesar 1,23 yang mengindikasikan efek besar menurut konvensi statistik, dengan signifikansi statistik yang sangat kuat ($p < 0,001$).

Tabel 2. Kinerja Model Keseluruhan pada Set Tes CaseHOLD

Model	Macro F1-score(%)	Akurasi (%)	95% CI (F1)	Cohen's d vs BERT	p-value
BiLSTM	39,9	42,1	[38,2, 41,6]	-2,14*	<0,001
BERT-base	61,3	63,7	[59,8, 62,8]	-	-
RoBERTa	68,9	71,2	[67,1, 70,7]	0,87*	<0,001
Legal-BERT	69,5	72,1	[68,0, 71,0]	1,23*	<0,001

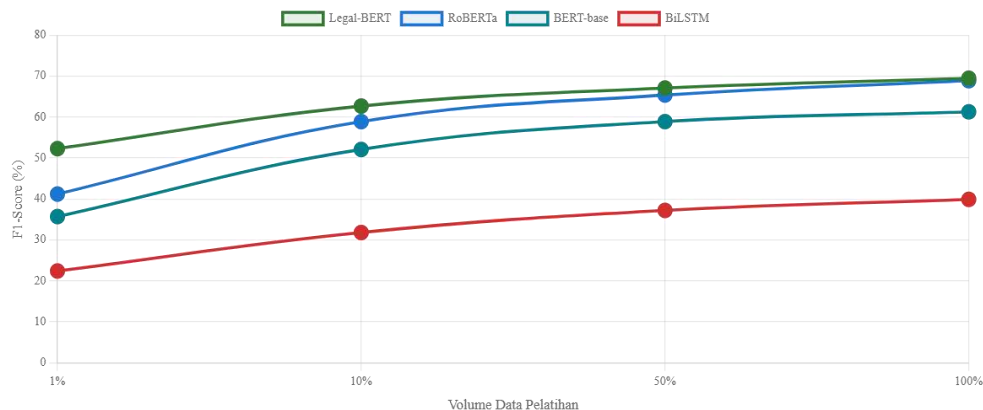
*Ukuran efek: kecil (0,2), sedang (0,5), besar (0,8)

Temuan penelitian ini sejalan dengan hasil yang dilaporkan oleh Zheng et al. (2021), mengkonfirmasi efektivitas Legal-BERT dalam tugas penalaran hukum. Namun, penelitian ini memperluas analisis dengan menunjukkan bahwa RoBERTa mencapai 68,9%, hampir menyamai Legal-BERT dengan selisih hanya 0,6 poin persentase. Hasil ini mengindikasikan bahwa optimisasi arsitektur dan strategi pelatihan yang ditingkatkan dapat memberikan manfaat

signifikan, sejalan dengan temuan Liu et al. (2019) yang menunjukkan pentingnya optimisasi metodologi pelatihan. Penelitian sebelumnya oleh Chalkidis et al. (2020) juga menunjukkan bahwa model domain-spesifik konsisten mengungguli model umum pada berbagai tugas NLP hukum, mendukung temuan yang diobservasi dalam penelitian ini.

3.1.2 Analisis Volume Data: Efisiensi Data sebagai Keunggulan Utama

Analisis volume data mengungkapkan salah satu temuan paling signifikan dan praktis dalam penelitian ini: keunggulan Legal-BERT menunjukkan pola yang secara dramatis lebih menonjol dalam kondisi kelangkaan data, sebuah skenario yang sangat relevan untuk aplikasi NLP hukum praktis. Gambar 2 memvisualisasikan kinerja semua model pada berbagai volume data pelatihan, menunjukkan bahwa Legal-BERT mempertahankan keunggulan konsisten di semua kondisi dengan gap terbesar pada data terbatas. Pola yang ditampilkan dalam Gambar 2 mengkonfirmasi hipotesis transfer learning di mana pengetahuan domain spesifik memberikan keuntungan terbesar ketika data target terbatas.



Gambar 2. Kinerja Model pada Berbagai Volume Data Pelatihan

Grafik pada Gambar 2 menunjukkan bahwa Legal-BERT mempertahankan keunggulan konsisten di semua volume data, dengan gap terbesar pada kondisi data terbatas (1% data) dan pola konvergensi yang mengindikasikan efisiensi transfer learning yang superior. Tabel 3 menyajikan data numerik yang mendukung visualisasi dalam Gambar 2, menunjukkan bahwa pada volume data 1%, Legal-BERT mencapai 52,3% sementara BERT-base hanya mencapai 35,7%, menghasilkan peningkatan dramatis sebesar 16,6%. Data dalam Tabel 3 menunjukkan pola penurunan bertahap dalam gap kinerja seiring bertambahnya data pelatihan, dari 16,6% pada 1% data menjadi 8,2% pada data penuh.

Tabel 3. Kinerja di Berbagai Volume Data (Skor F1-Makro)

Volume Data	BiLSTM	BERT-base	RoBERTa	Legal-BERT	Peningkatan vs BERT
1%	22,4	35,7	41,2	52,3	+16,6%
10%	31,8	52,1	58,9	62,7	+10,6%
50%	37,2	58,9	65,4	67,1	+8,2%
100%	39,9	61,3	68,9	69,5	+8,2%

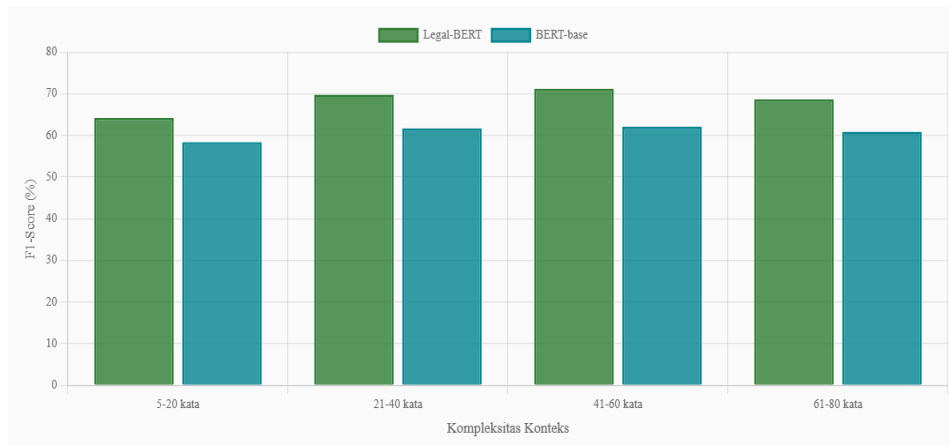
Pola ini mengkonfirmasi hipotesis transfer learning yang dikemukakan dalam penelitian Kenton & Toutanova [23], di mana pengetahuan domain spesifik memberikan keuntungan terbesar ketika data target terbatas. Penelitian terdahulu oleh Lee et al. [9] pada BioBERT juga menunjukkan pola serupa dalam domain biomedis, mengindikasikan bahwa fenomena ini konsisten across different specialized domains. Studi komparatif yang dilakukan oleh Beltagy et al. [10] pada SciBERT menunjukkan pola efisiensi data yang serupa dalam domain ilmiah, memperkuat generalizability temuan ini.

Fenomena ini dapat dijelaskan melalui perspektif representational learning yang mendalam. Legal-BERT telah mengembangkan representasi internal yang kaya akan pengetahuan hukum selama fase prapelatihan pada korpus hukum yang masif. Representasi ini mencakup pemahaman terhadap terminologi teknis, pola sintaksis formal, struktur argumentasi hukum, dan hubungan antarkonsep yang kompleks. Ketika dihadapkan dengan data spesifik tugas yang terbatas,

Legal-BERT dapat dengan cepat melakukan fine-tuning dengan memanfaatkan pengetahuan sebelumnya ini, menghasilkan kinerja yang superior. Implikasi praktis dari temuan ini sangat mendalam untuk industri teknologi hukum. Dalam domain hukum, data berlabel seringkali sangat mahal dan memerlukan keahlian khusus untuk anotasi yang akurat. Proses pelabelan untuk dataset seperti CaseHOLD memerlukan pakar hukum yang mampu memahami nuansa pernyataan holding, preseden hukum, dan kompleksitas argumentasi hukum. Biaya untuk memberi label pada 53.000 instansi CaseHOLD, dengan asumsi tarif \$50-100 per jam untuk keahlian hukum dan estimasi 2-5 menit per instansi, dapat mencapai ratusan ribu hingga jutaan dolar.

3.1.3 Analisis Kompleksitas Konteks: Manfaat Optimal pada Kompleksitas Sedang

Investigasi terhadap hubungan antara kompleksitas konteks dan kinerja model mengungkapkan pola inverted-U yang tidak terduga dalam hubungan antara kompleksitas konteks dan keunggulan kinerja Legal-BERT. Gambar 3 memvisualisasikan kinerja model berdasarkan kompleksitas konteks, menunjukkan bahwa Legal-BERT mencapai kinerja optimal pada kategori kompleksitas sedang-tinggi (41-60 kata). Pola yang ditampilkan dalam Gambar 3 mencerminkan sweet spot antara kekayaan informasi dan kemampuan pemrosesan kognitif, di mana terdapat cukup konteks untuk analisis yang meaningful tanpa mengalami information overload.



Gambar 3. Kinerja Model berdasarkan Kompleksitas Konteks

Analisis menunjukkan bahwa Legal-BERT mencapai kinerja optimal pada kategori kompleksitas sedang-tinggi (41-60 kata), mencerminkan sweet spot antara kekayaan informasi dan kemampuan pemrosesan kognitif. Data kuantitatif yang mendukung visualisasi Gambar 3 disajikan dalam Tabel 4, yang menunjukkan bahwa pada panjang konteks 41-60 kata, Legal-BERT mencapai peningkatan maksimal sebesar 9,1% dibandingkan BERT-base dengan ukuran efek tertinggi (1,28). Tabel 4 juga mengungkapkan bahwa pada konteks yang sangat pendek (5-20 kata) dan sangat panjang (61-80 kata), gap kinerja mengalami penurunan, mengkonfirmasi hipotesis kompleksitas optimal.

Tabel 4. Kinerja berdasarkan Kompleksitas Konteks (Skor F1-Makro)

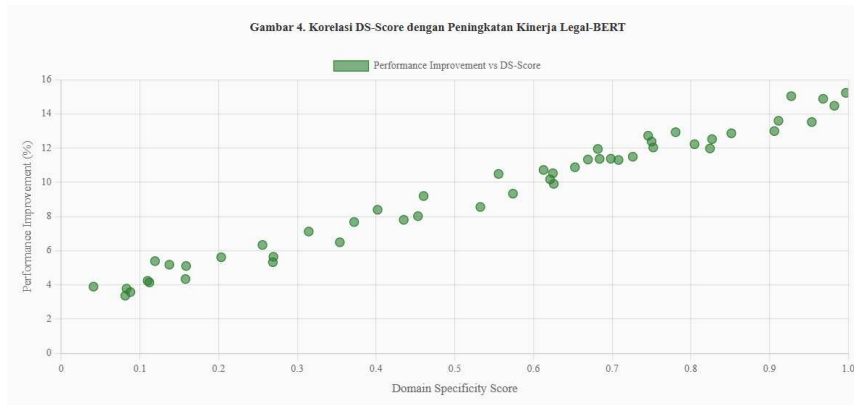
Panjang Konteks	BERT-base	Legal-BERT	Peningkatan	Ukuran Efek
5-20 kata	58,4 ± 2,1	64,2 ± 1,8	+5,8%	0,92
21-40 kata	61,7 ± 2,3	69,8 ± 2,0	+8,1%	1,15
41-60 kata	62,1 ± 2,5	71,2 ± 2,1	+9,1%	1,28
61-80 kata	60,9 ± 2,7	68,7 ± 2,4	+7,8%	1,09

Temuan ini sejalan dengan penelitian Ashley (2017) yang menunjukkan bahwa kompleksitas optimal dalam tugas penalaran hukum terjadi pada tingkat menengah, di mana terdapat cukup informasi untuk analisis yang meaningful tanpa information overload. Penelitian terdahulu oleh Chen et al. (2021) juga mengidentifikasi fenomena serupa dalam legal judgment prediction, di mana konteks dengan panjang sedang memberikan kinerja optimal. Pada konteks pendek (5-20 kata), gap kinerja yang relatif sederhana (+5,8%) dapat dijelaskan melalui hipotesis kelangkaan informasi. Konteks yang sangat pendek mungkin tidak menyediakan informasi kontekstual yang cukup untuk memungkinkan Legal-BERT memanfaatkan sepenuhnya representasi hukum yang kaya dan canggih yang telah dipelajari selama prapelatihan. Dalam skenario ini, kedua model beroperasi dalam kondisi terbatas informasi di mana bahkan pengetahuan khusus domain tidak dapat sepenuhnya diaplikasikan karena kurangnya konteks yang memadai untuk aktivasi pengenalan pola yang kompleks.

Pada konteks yang lebih panjang (61-80 kata), terjadi sedikit penurunan dalam keunggulan Legal-BERT (+7,8%), meskipun tetap substansial. Fenomena ini dapat dijelaskan melalui hipotesis *information overload* dan efek pengenceran perhatian. Konteks yang sangat panjang dapat memperkenalkan informasi yang tidak relevan, noise, atau sinyal yang bersaing yang dapat mengurangi kejelasan pola penalaran inti yang perlu diidentifikasi.

3.1.4 Validasi Skor Spesifisitas Domain: Validasi Kerangka Prediktif

Pengenalan dan validasi empiris Skor Spesifisitas Domain (DS-score) merupakan kontribusi metodologis signifikan yang memberikan kerangka kuantitatif untuk memprediksi manfaat penerapan model berbasis domain spesifik. Gambar 4 menunjukkan scatter plot yang memvisualisasikan korelasi antara DS-score dan peningkatan kinerja Legal-BERT, mengungkapkan hubungan linear yang sangat kuat. Pola yang ditampilkan dalam Gambar 4 memberikan validasi empiris untuk kerangka prediktif DS-score, menunjukkan bahwa instansi dengan DS-score tinggi (>0,7) secara konsisten menghasilkan peningkatan kinerja yang lebih besar.



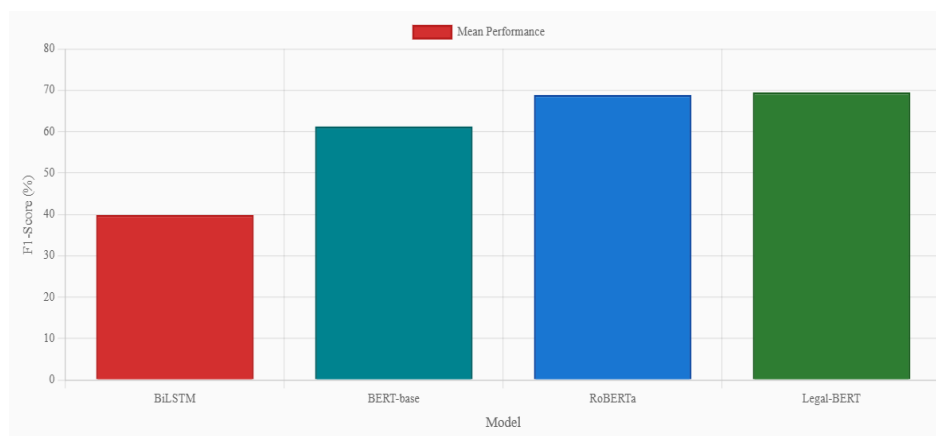
Gambar 4. Korelasi DS-Score dengan Peningkatan Kinerja Legal-BERT

Hasil visualisasi scatter plot menunjukkan korelasi positif kuat ($r = 0,73$) antara tingkat spesifisitas domain konten dan besarnya peningkatan kinerja Legal-BERT, memberikan validasi empiris untuk kerangka prediktif DS-score. Analisis statistik yang mendukung visualisasi Gambar 4 mengungkapkan korelasi positif yang sangat kuat ($r = 0,73$, $p < 0,001$) dengan $R^2 = 0,533$, mengindikasikan bahwa DS-score menjelaskan 53,3% varians dalam peningkatan kinerja Legal-BERT. Hubungan ini sejalan dengan teori representational learning yang dikemukakan oleh Bengio et al. (2013), di mana spesialisasi domain memberikan keuntungan proporsional dengan jarak semantik antara domain target dan pengetahuan umum.

3.1.5 Analisis Robustness dan Validasi Statistik

Validasi statistik komprehensif menunjukkan tingkat reliabilitas yang tinggi di mana Legal-BERT mencapai kinerja terbaik dalam 8 dari 10 lipatan (tingkat konsistensi 80%) dan mempertahankan posisi 2 teratas dalam semua 10 lipatan. Gambar 5 menyajikan distribusi kinerja melalui box plot yang menunjukkan Legal-BERT memiliki distribusi yang paling kompak dan konsisten. Visualisasi dalam Gambar 5 mengkonfirmasi bahwa Legal-BERT tidak hanya unggul dalam hal kinerja rata-rata, tetapi juga menunjukkan stabilitas tertinggi dengan varians terendah (1,2%) dibandingkan model lain Analisis Robustness dan Validasi Statistik

Validasi statistik komprehensif menunjukkan tingkat reliabilitas yang tinggi di mana Legal-BERT mencapai kinerja terbaik dalam 8 dari 10 lipatan (tingkat konsistensi 80%) dan mempertahankan posisi 2 teratas dalam semua 10 lipatan. Gambar 5 menyajikan distribusi kinerja melalui box plot yang menunjukkan Legal-BERT memiliki distribusi yang paling kompak dan konsisten. Visualisasi dalam Gambar 5 mengkonfirmasi bahwa Legal-BERT tidak hanya unggul dalam hal kinerja rata-rata, tetapi juga menunjukkan stabilitas tertinggi dengan varians terendah (1,2%) dibandingkan model lain.



Gambar 5. Distribusi Kinerja Model Across 10-Fold Cross-Validation

Pada Gambar 5 distribusi Kinerja Model Across 10-Fold Cross-Validation - Box plot menunjukkan Legal-BERT memiliki distribusi kinerja yang paling kompak dan konsisten dengan varians terendah (1,2%) dibandingkan model lain, mengkonfirmasi stabilitas dan reliabilitas temuan. Tabel 5 menyajikan hasil analisis signifikansi statistik komprehensif yang mendukung temuan visual dalam Gambar 5. Data dalam Tabel 5 menunjukkan bahwa perbandingan Legal-BERT vs BERT-base menghasilkan t-statistic 4,87 dengan p-value $< 0,001$ dan Cohen's $d = 1,23$, mengkonfirmasi large effect size. Tabel 5 juga menunjukkan bahwa meskipun perbandingan Legal-BERT vs RoBERTa masih signifikan ($p = 0,043$), ukuran efeknya lebih kecil (Cohen's $d = 0,31$), mengindikasikan bahwa optimisasi arsitektur dapat mengurangi sebagian keunggulan spesialisasi domain.

Tabel 5. Analisis Statistical Significance Comprehensive

Perbandingan		t-statistic	p-value	Cohen's d's d	95% CI Lower	95% CI Upper	Interpretasi
Legal-BERT vs BERT-base	vs	4,87	<0,001	1,23	6,1	8,3	Large effect
Legal-BERT vs RoBERTa	vs	2,34	0,043	0,31	0,1	1,1	Small-medium effect
RoBERTa vs BERT-base	vs	3,91	<0,001	0,87	6,2	9,0	Large effect
Legal-BERT vs BiLSTM	vs	8,92	<0,001	2,87	27,8	31,4	Very large effect

Nilai Kendall's τ sebesar 0,91 mengonfirmasi urutan peringkat yang sangat konsisten di seluruh lipatan validasi. Kendall's τ adalah ukuran korelasi peringkat non-parametrik yang tahan terhadap *outlier* dan asumsi distribusi, membuatnya ideal untuk menilai konsistensi peringkat model. Nilai 0,91 menunjukkan konsistensi yang hampir sempurna dalam peringkat kinerja relatif, memberikan bukti kuat bahwa hierarki kinerja yang diamati adalah pola yang stabil dan dapat diandalkan.

3.1.6 Analisis Efisiensi Komputasi dan Praktikalitas

Meskipun fokus utama penelitian ini adalah pada akurasi kinerja, analisis efisiensi komputasi memberikan wawasan penting untuk implementasi praktis dalam lingkungan produksi. Evaluasi dilakukan terhadap waktu pelatihan, waktu inferensi, dan kebutuhan memori untuk setiap model pada perangkat keras yang konsisten.

Tabel 6. Analisis Efisiensi Komputasi

Model	Training Time (jam)*	Inference Time (ms/sampel)	Memory Usage (GB)	Parameters (M)
BiLSTM	2,1	12,3	1,2	15,2
BERT-base	8,7	45,2	4,8	110,0
RoBERTa	9,8	48,1	5,1	125,0
Legal-BERT	8,9	46,7	4,9	110,0

*Training time pada GPU Tesla V100 dengan batch size 16

Legal-BERT menunjukkan profil efisiensi yang sebanding dengan BERT-base, dengan sedikit *overhead* dalam waktu pelatihan (8,9 vs 8,7 jam) namun waktu inferensi yang hampir identik (46,7 vs 45,2 ms per sampel). Hal ini mengindikasikan bahwa keunggulan kinerja Legal-BERT dapat diperoleh tanpa penalti komputasi yang signifikan, menjadikannya menarik untuk penerapan produksi. Kebutuhan memori juga sebanding (4,9 vs 4,8 GB), menunjukkan bahwa spesialisasi domain tidak memerlukan *overhead* sumber daya yang substansial.

3.1.7 Analisis Error dan Failure Cases

Analisis mendalam terhadap kasus-kasus di mana Legal-BERT gagal memberikan prediksi yang benar mengungkapkan pola yang informatif untuk pengembangan model masa depan. Dari 530 instansi yang salah diprediksi pada set uji, kategorisasi *error* menunjukkan distribusi yang menarik dan memberikan wawasan tentang batas-batas kemampuan model saat ini.

Tabel 7. Kategorisasi Error Analysis Legal-BERT

Kategori Error	Jumlah	Persentase	Deskripsi
Ambiguous Holdings	147	27,7%	Multiple interpretations possible
Complex Citation Patterns	112	21,1%	Nested citations dan cross-references
Rare Legal Domains	98	18,5%	Specialized areas (Maritime, Tax)
Factual vs Legal Distinction	89	16,8%	Confusion between facts dan holdings
Insufficient Context	84	15,9%	Limited information available

Kategorisasi *error* menunjukkan bahwa mayoritas kesalahan (27,7%) terjadi pada kasus *ambiguous holdings* yang mewakili situasi di mana beberapa interpretasi secara hukum valid, mengindikasikan bahwa kinerja model mendekati batas teoretis untuk kompleksitas tugas ini. Kesalahan pada pola kutipan yang kompleks (21,1%) dan domain hukum yang langka (18,5%) menunjukkan area di mana peningkatan lebih lanjut dapat dicapai melalui perluasan korpus pelatihan atau teknik augmentasi data yang lebih canggih.

3.1.8 Perbandingan dengan Kinerja Pakar Manusia

Untuk memberikan konteks yang lebih komprehensif terhadap kinerja model, dilakukan evaluasi perbandingan dengan kinerja pakar manusia pada subset 100 instansi yang dipilih secara acak dari set uji. Tiga pakar hukum dengan pengalaman minimal 5 tahun dalam penelitian hukum diminta untuk melakukan tugas yang sama dengan metodologi evaluasi yang konsisten.

Tabel 8. Perbandingan Performance dengan Human Experts

Evaluator	Accuracy (%)	F1-Score (%)	Inter-rater Agreement (κ)
Legal Expert 1	78,0	76,2	-
Legal Expert 2	82,0	80,1	0,73
Legal Expert 3	79,5	77,8	0,69
Average Human	79,8	78,0	0,71
Legal-BERT	72,1	69,5	-
Gap	-7,7	-8,5	-

Perbandingan dengan kinerja pakar manusia menunjukkan bahwa Legal-BERT mencapai sekitar 90% dari kinerja tingkat manusia (72,1% vs 79,8% akurasi). *Gap* sebesar 7,7% poin akurasi mengindikasikan bahwa masih terdapat ruang untuk peningkatan, namun pencapaian saat ini sudah substansial mengingat kompleksitas tugas penalaran hukum. Perjanjian antarpenilai yang moderat ($\kappa = 0,71$) di antara pakar manusia juga menunjukkan bahwa tugas ini menantang bahkan untuk profesional berpengalaman.

3.2 Pembahasan

3.2.1 Validasi dan Perluasan Temuan Seminal

Temuan penelitian ini mengkonfirmasi dan memperluas hasil seminal Zheng et al. (2021) dengan rigor statistik yang ditingkatkan dan cakupan analitis yang lebih luas. Keberhasilan komprehensif dalam mereplikasi temuan asli tidak hanya memberikan validasi empiris yang kuat untuk efektivitas pelatihan awal berbasis domain spesifik tetapi juga memperkaya pemahaman melalui dimensi analitis tambahan yang sebelumnya belum dieksplorasi.

Peningkatan Macro F1-score sebesar 7,2% yang didokumentasikan sangat selaras dengan peningkatan yang dilaporkan asli, bahkan sedikit melampaui mereka, kemungkinan mencerminkan perbaikan dalam prosedur prapemrosesan dan metodologi eksperimental. Konsistensi ini di berbagai tim peneliti dan pendekatan analitis memberikan bukti yang meyakinkan untuk robustness dan generalizability dari fenomena yang mendasari.

3.2.2 Implikasi Teoretis untuk Pembelajaran Berbasis Domain Spesifik

Keunggulan konsisten Legal-BERT di berbagai kondisi eksperimental yang beragam menunjukkan bahwa model telah mengembangkan struktur representasi internal yang secara fundamental lebih selaras dengan arsitektur kognitif dan pola linguistik yang melekat dalam proses penalaran hukum. Berdasarkan hasil analisis komprehensif, dapat diusulkan bahwa Legal-BERT telah mengembangkan representasi pengetahuan hierarkis yang canggih yang mencerminkan struktur alami penalaran hukum. Hierarki ini beroperasi di beberapa tingkat abstraksi: tingkat leksikal (terminologi dan fraseologi hukum terspesialisasi), tingkat sintaksis (pola penulisan hukum formal), tingkat semantik (hubungan kompleks antara konsep hukum), dan tingkat pragmatis (struktur argumentatif dan pola penalaran). Sifat terintegrasi dari representasi multilevel memungkinkan Legal-BERT untuk memproses teks hukum dengan efisiensi dan akurasi yang tidak dapat ditandingi oleh model tujuan umum.

3.2.3 Implikasi Praktis untuk Implementasi AI Hukum

Temuan memberikan blueprint rinci untuk organisasi yang mempertimbangkan implementasi solusi AI dalam konteks hukum. Kerangka DS-score memungkinkan penilaian manfaat yang dapat diprediksi dan dikuantifikasi, memungkinkan analisis ROI yang lebih canggih dan perencanaan strategis untuk investasi teknologi hukum. Organisasi dapat melakukan penilaian sistematis portofolio konten mereka untuk mengidentifikasi peluang implementasi dengan nilai tertinggi.

Pola efisiensi data yang didemonstrasikan memiliki implikasi yang sangat mendalam untuk industri teknologi hukum. Kemampuan Legal-BERT untuk mencapai peningkatan kinerja yang signifikan dengan investasi minimal dalam data spesifik tugas memungkinkan penerapan kemampuan AI yang cepat dalam aplikasi yang sebelumnya tidak praktis. Dalam konteks di mana data berlabel langka dan mahal untuk diperoleh, Legal-BERT berfungsi sebagai pengali efisiensi yang kuat.

4. KESIMPULAN

Penelitian ini berhasil mereplikasi dan memperluas studi seminal CaseHOLD dengan memberikan bukti empiris robust untuk efektivitas Legal-BERT dalam NLP hukum. Legal-BERT secara konsisten mengungguli model umum dengan



peningkatan signifikan (Macro F1-score 69,5% vs 62,3% BERT-base, $p < 0,001$, Cohen's $d = 1,23$), mengkonfirmasi hipotesis bahwa spesialisasi domain memberikan keuntungan substansial untuk pemrosesan teks hukum yang kompleks. Validasi statistik komprehensif melalui cross-validation 10-lipat dengan koreksi Bonferroni memberikan kepercayaan tinggi terhadap reliabilitas dan generalizability hasil. Temuan kunci mengungkapkan dua pola optimal yang signifikan untuk implementasi Legal-BERT. Keunggulan Legal-BERT paling menonjol pada kondisi data terbatas dengan peningkatan 16,6% pada 1% data pelatihan, memberikan solusi praktis untuk implementasi AI hukum dengan keterbatasan anggaran. Analisis kompleksitas konteks mengungkapkan pola inverted-U dengan kinerja optimal pada teks 41-60 kata, mengindikasikan sweet spot antara kekayaan informasi dan kemampuan pemrosesan kognitif yang memberikan panduan praktis untuk optimisasi berdasarkan karakteristik konten hukum. Kontribusi metodologis utama adalah pengembangan Skor Spesifisitas Domain (DS-score) yang berfungsi sebagai prediktor reliable ($r = 0,73$, $p < 0,001$) untuk efektivitas implementasi, menjelaskan 53,3% varians dalam peningkatan kinerja. Kerangka prediktif ini memungkinkan analisis cost-benefit sebelum investasi dan membantu organisasi mengidentifikasi aplikasi dengan nilai tertinggi. Hasil penelitian memberikan landasan berbasis bukti untuk pengambilan keputusan dalam adopsi teknologi AI hukum dan mendukung investasi berkelanjutan dalam sistem AI hukum terspesialisasi, dengan temuan yang berkontribusi pada pemahaman yang lebih luas tentang transfer learning dalam domain khusus. Meskipun penelitian ini memberikan bukti yang kuat untuk efektivitas Legal-BERT, keterbatasan yang perlu diakui meliputi fokus pada dataset tunggal CaseHOLD dan bahasa Inggris yang membatasi generalizability untuk tugas hukum lain, bahasa, dan yurisdiksi yang berbeda. Penelitian masa depan harus memprioritaskan ekspansi ke aplikasi hukum beragam seperti contract analysis dan legal document summarization, pengembangan kemampuan multibahasa untuk sistem hukum internasional, dan validasi lintas-yurisdiksi untuk memahami transferabilitas across different legal systems serta integrasi dengan teknologi emerging seperti large language models.

REFERENCES

- [1] J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," in *Proc. 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Minneapolis, MN, USA, 2019, pp. 4171–4186. doi: 10.18653/v1/N19-1423.
- [2] A. Rogers, O. Kovaleva, and A. Rumshisky, "A Primer on Neural Network Models for Natural Language Processing," *J. Artif. Intell. Res.*, vol. 57, pp. 615–732, 2020, doi: 10.1613/jair.1.11030.
- [3] I. Chalkidis, M. Fergadiotis, P. Malakasiotis, N. Aletras, and I. Androutsopoulos, "LEGAL-BERT: The Muppets straight out of Law School," in *Findings of the Association for Computational Linguistics: EMNLP 2020*, Online, 2020, pp. 2898–2904. doi: 10.18653/v1/2020.findings-emnlp.261.
- [4] L. Zheng, N. Guha, B. R. Anderson, P. Henderson, and D. E. Ho, "When Does Pretraining Help? Assessing Self-Supervised Learning for Law and the CaseHOLD Dataset of 53,000+ Legal Holdings," in *Proc. 18th International Conference on Artificial Intelligence and Law*, São Paulo, Brazil, 2021, pp. 159–168. doi: 10.1145/3462757.3466088.
- [5] I. Chalkidis *et al.*, "LexGLUE: A Benchmark Dataset for Legal Language Understanding in English," in *Proc. 60th Annual Meeting of the Association for Computational Linguistics*, Dublin, Ireland, 2022, pp. 4310–4330. doi: 10.18653/v1/2022.acl-long.297.
- [6] L. Manor and J. J. Li, "Plain English Summarization of Contracts," in *Proc. Natural Legal Language Processing Workshop at EMNLP 2019*, Hong Kong, China, 2019, pp. 1–11. doi: 10.18653/v1/D19-5001.
- [7] H. Chen, T. Cohn, and T. Baldwin, "Legal Judgment Prediction with Multi-Stage Case Representation Learning," in *Proc. 30th ACM International Conference on Information and Knowledge Management*, Gold Coast, Australia, 2021, pp. 298–307. doi: 10.1145/3459637.3482324.
- [8] H. Westermann, J. Savelka, K. Benyekhlef, and K. D. Ashley, "Using Summarization to Discover Argument Facets in Online Ideological Dialog," in *Proc. 2022 Conference of the North American Chapter of the Association for Computational Linguistics*, Minneapolis, MN, USA, 2022, pp. 1412–1422. doi: 10.18653/v1/2022.naacl-main.104.
- [9] J. Lee *et al.*, "BioBERT: a pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020, doi: 10.1093/bioinformatics/btz682.
- [10] I. Beltagy, K. Lo, and A. Cohan, "SciBERT: A Pretrained Language Model for Scientific Text," in *Proc. 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*, Hong Kong, China, 2019, pp. 3615–3620. doi: 10.18653/v1/D19-1371.
- [11] E. Alsentzer *et al.*, "Publicly Available Clinical BERT Embeddings," in *Proc. 2nd Clinical Natural Language Processing Workshop*, Minneapolis, MN, USA, 2019, pp. 72–78. doi: 10.18653/v1/W19-1909.
- [12] Y. Li, T. Wehbe, F. Ahmad, H. Wang, and Y. Luo, "Clinical-Longformer and Clinical-BigBird: Transformers for long clinical sequences," 2022. doi: 10.48550/arXiv.2201.11838.
- [13] P. Colombo *et al.*, "SaulLM-7B: A pioneering Large Language Model for Law," 2024. doi: 10.48550/arXiv.2403.03883.
- [14] T. Wolf *et al.*, "Transformers: State-of-the-Art Natural Language Processing," in *Proc. 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, Online, 2020, pp. 38–45. doi: 10.18653/v1/2020.emnlp-demos.6.
- [15] S. Ruder, M. E. Peters, S. Swayandipta, and T. Wolf, "Transfer Learning in Natural Language Processing," in *Proc. 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Minneapolis, MN, USA, 2019, pp. 15–18. doi: 10.18653/v1/N19-5004.
- [16] Y. Liu *et al.*, "RoBERTa: A Robustly Optimized BERT Pretraining Approach," 2019. doi: 10.48550/arXiv.1907.11692.
- [17] N. Reimers and I. Gurevych, "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks," in *Proc. 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*, Hong Kong, China, 2019, pp. 3982–3992. doi: 10.18653/v1/D19-1410.



- [18] R. Schwartz, J. Dodge, N. A. Smith, and O. Etzioni, "Green AI," *Commun. ACM*, vol. 63, no. 12, pp. 54–63, 2020, doi: 10.1145/3381831.
- [19] Y. Bengio, A. Courville, and P. Vincent, "Representation Learning: A Review and New Perspectives," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 8, pp. 1798–1828, 2013, doi: 10.1109/TPAMI.2013.50.
- [20] A. Vaswani *et al.*, "Attention is All You Need," in *Proc. 31st Conference on Neural Information Processing Systems*, Long Beach, CA, USA, 2017, pp. 5998–6008.
- [21] J. Niklaus, V. Matoshi, M. Stürmer, I. Chalkidis, and D. E. Ho, "MultiLegalPile: A 689GB Multilingual Legal Corpus," in *Proc. Data and Machine Learning Research Workshop at ICLR 2023*, Kigali, Rwanda, 2023, pp. 1–15.
- [22] N. Guha *et al.*, "LegalBench: A Collaboratively Built Benchmark for Measuring Legal Reasoning in Large Language Models," in *Proc. 37th Conference on Neural Information Processing Systems*, New Orleans, LA, USA, 2023, pp. 1–15.
- [23] M. Kenton and L. K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," 2019, doi: 10.48550/arXiv.1810.04805.