

スマートフォンでのニュース記事閲覧における LLM 要約長の効果：主観評価と理解度に基づく最適文字量の検討

鈴木 健司¹ 坂本 大介²

概要：本研究は、スマートフォンでのニュース記事閲覧を対象に、大規模言語モデル (LLM) が生成する要約の文字数が、読者の主観印象および理解に及ぼす影響を検討する。LLM 要約は一般化しつつある一方、スマートフォンの表示制約を踏まえた適切な要約文字数に関する体系的知見は十分でない。そこで、要約文字数を 7 水準で操作し、スマートフォン上で提示した要約に対して、主観評価（読みやすさ、情報の充実度、面白さ、読疲れ感）および理解度テストを実施した。その結果、「情報の充実度」は 300~400 字付近で飽和する傾向が見られ、「面白さ」は文字数の増加に伴い上昇した。一方で、「読疲れ感」も文字数に応じて増加し、特に 800 字を超えると上昇が顕著となった。理解度テストでは要約文字数による有意差は認められなかった。以上より、スマートフォン向けニュース記事要約では、情報量と読解負荷のバランスが比較的良好となる文字数は概ね 400 字付近であることが示唆される。本知見は、ニュース配信サービスや自動要約システムにおける要約文字数設計の指針として有用である。

1. はじめに

スマートフォンの普及により、日常的なニュース閲覧や情報収集の多くが手軽に持ち歩けるスマートフォンのタッチディスプレイ上で行われるようになった。このようなモバイル環境では、比較的小型の限られた表示領域をもつディスプレイと、移動に伴う閲覧時間の制約から、短時間で要点を把握する需要が高まりつつある [1, 2]。このような状況に対し、これまで要約は人間の手作業で行われてきたが、大規模言語モデル (Large Language Model; LLM) の文章の自動要約技術の発展により、自動で長いニュース記事やレポートなどの文章を要約し、文章量を短縮することが可能となってきた。

LLM による文章の要約は、ニュース記事やレポートの著者や編集者の作業効率化に大きく寄与するものであると考えられる一方で、LLM による文章の要約が文章を読む体験にどのような影響を与えるかについては十分に調査がされてきていない。例えば、LLM による文章の要約が短すぎる場合には情報の欠落により理解が困難となる一方、長すぎる場合には読解負荷が高まり、読了率や満足度が低下する可能性がある。特に、スマートフォンのように画面サイズが小さい環境では、スクロール量や視線移動といった物理的制約が読者の心理的負担に大きく影響することが

知られている [3-5]。したがって、スマートフォン環境における要約文章の最適な文字量を明らかにすることは、情報提示の設計や自動要約システムの実用化において重要な課題である。

本研究では、ニュース記事を題材とし、LLM によって異なる文字数に要約された文章を用いて、スマートフォン上での主観評価および理解度の比較実験を行う。要約文字数を 100~1000 字の範囲で 7 条件設定し、クラウドソーシングを用いた大規模調査を実施する。スクリーニング後の有効参加者数は 2709 名である。調査によって得られた読みやすさ、理解しやすさ、情報の充実度、面白さ、疲労感といった主観的印象、および内容理解の程度を定量的に分析することで、スマートフォンでのニュース記事閲覧時の LLM による文章要約の効果を検討する。結果として、LLM による文章要約時の適切な文字量とその効果について検討する。本論文の最後では、本研究で得られる知見をニュース配信を含む要約重視型サービスや自動要約システムの設計指針として利活用するための方針について議論する。

2. 関連研究

スマートフォン環境における文章量に関する研究は、これまでに様々な研究が行われている。

2.1 小型ディスプレイにおける読書

画面サイズや行長、行送りなどの物理的条件は、スクリー

¹ LINE ヤフー株式会社

² 北海道大学

ンにおける読書の速度や理解、疲労に影響することが古くから指摘されている [6, 7]. 小型ディスプレイ環境に焦点を当てた研究では、スクロールやページ送りなどのユーザ操作は理解や記憶、視覚疲労に影響を与えると報告されている [3-5]. 特に小さいディスプレイでの縦スクロールは、読解時の戻り読みの行い方や統合的理解に影響を与えることが指摘されている [5]. またスマートフォン環境での読書は、紙面と異なる生理反応の負荷特性（ため息発生の抑制、前頭前野活動の亢進等）を示し、理解低下に関与しうることが示唆されている [8]. 一方で、フォントサイズや行間の調整によって視線指標や主観評価は改善される可能性がある [9].

2.2 モバイル環境におけるニュース閲覧と要約

ニュースを閲覧する環境はこの数年で更にモバイル・ソーシャルにシフトしており [1, 2], 小画面・短時間での要点把握の重要性が増している. こうした中、画面サイズに応じて要点を段階的に提示するレスポンス要約やスキムリーディング支援の UI の提案や [10], ページングによる分割表示とスクロールによる全文表示の選択が読書体験に与える影響について再検討されている [4, 5]. 本研究は、こうした UI 要素における「要約文章の分量そのもの」に焦点を絞り、同一記事の要約を異なる文字数で提示した時に、主観評価と理解度がどう変わるかを大規模に検証する点に目的がある.

2.3 LLM による要約の品質と評価

LLM を含むニューラル要約は高い流暢性を示す一方、事実性や一貫性に課題があることが広く報告されており [11], 自動評価尺度の限界と人手評価の設計が議論されてきた [12]. 事実整合性の自動評価では、生成された要約中の事実上の不整合を特定する「QAGS」などが提案されたが [13], 落ちた情報の網羅性は直接評価しない. 本研究は、スマートフォンという制約下で要約の長さが主観評価（充実度、面白さ、疲労感）と理解度（主要な考えや事実）に与える影響を総合的に測定する.

実運用ではメディアや UI デザインの都合で所定の長さに要約文章量を合わせる必要がある. ニューラル要約での文章量制御は、デコーディング時の制約付与と学習時の長さ埋め込みや位置符号化拡張に大別される [14, 15]. さらに属性制御により、要約の文章量やスタイル、焦点を当てる対象、あるいは未読部分のみを指定して生成する手法が提案されている. [16]. 本研究の設定（100~1000 字の広範囲での主観・理解度の比較）は、こうした制御手法に対してどの文章量が望ましいかという UI デザイン観点の外部知見を提供する.

2.4 関連研究のまとめと本研究の目的

総じて、小型ディスプレイでは文章の長さやレイアウト、遷移様式が相互作用的に読書体験を規定しうる [5, 6, 9]. 一方、ニューラル要約は、文章の長さや事実性・可読性のトレードオフがあり [11, 12], 適切な文章の長さに関する知見は多くない. 本研究はニュース記事を対象に、スマートフォンという現実的な制約下で要約文章の文章量による主観評価と理解度の関係を大規模に検証し、実運用のための目安となる文章量を示すことを目的とする.

3. LLM による文章要約と要約された文章の詳細

本研究においては主にニュースサイトの記事を開覧する状況を扱うため、実際のニュースサイトの記事を要約する文章として使用することとした. 当研究グループの事前調査では、公開されているニュースサイトの多くは記事本文が 1000 字程度であることを確認した. 実験で使用される記事の文章を要約する文字数について検討するために、本研究ではまず 100, 200, 300, 400 字を対象に読みやすさに関する主観評価と文章の主旨に関する理解度テストを実施した. その結果、主観評価において 400 字以上に最適値がある可能性が示唆されたため、400 字を超える文字数として 600, 800, 1000 字についての追加の調査を行った. 400 字までは 100 字間隔では大きな差が見られなかったため、追加の調査では 200 字間隔とした. この結果、元々の記事の長さである 1000 字を含めた 100, 200, 300, 400, 600, 800, 1000 字の 7 種類の要約文章を用意することとした.



図 1: 実験で用いた要約文章の表示例
(左から 200 字, 400 字, 600 字)

権利関係で実際に使用した文章を本稿で使用できないため、ここでは Wikipedia の文章を使用している（脚注参照）.

実験で用いた要約文章の表示例を図 1 に示す. 調査で実際に使用したニュース記事は権利の関係で本稿では使用できないため、ここでは Wikipedia の文章を使用した*1. 要約文章を表示する際のフォントプロパティは、Yahoo!

*1 <https://ja.wikipedia.org/wiki/可読性>

JAPAN が運営する Yahoo!ニュースの設定と同様とした。すなわち、フォントサイズは 17px、ラインハイトは 1.5、余白は上左右で 8px とした。Yahoo!ニュースは月間 225 億 PV のアクセスがある日本で最も利用されているニュースサービスの一つであり [17]、検証を行うベースラインとして適切だと考えた。本実験の独自 UI である読了完了ボタンは、要約文章がディスプレイ高を超えない場合はディスプレイ下部に配置し、超えた場合は要約文章下部から 1em の余白を空けた。

3.1 LLM による文章要約の方法

本実験には Yahoo!ニュースに掲載された日本語で書かれた記事を用いる。記憶による影響を考慮し、掲載後 1 ヶ月以上経過した記事を対象とした。実験に用いる文章の内容について、Dyson らは、特定の専門分野に偏らず「一般的な関心」を引く内容の文章を用意することで、実験変数への影響を正確に測定しようとした [18]。Rello らは、技術的・学術的すぎず、一般的な関心を持たれやすいトピックスを選んだ [9]。本実験ではこれらの知見を踏まえ、記事内容に対する個人的賛否や嗜好が主観評価に影響することを抑制するため、賛否が分かれにくく対立性が低い題材として、地域の話題やクリエイターへのインタビュー記事から選抜した。さらに、特定の記事内容に起因する偶発的な効果が結果に過度な影響を与えることを防ぐため、単一の記事ではなく複数の記事を用いる実験設計とし、9 本の記事を用意した。本実験では、これらの記事から生成された要約文章に対し主観評価と理解度テストを行うことから、元記事間の難易度は揃っていることが望ましい。そこで佐藤の日本語文章難易度測定 [19] を行った。この測定手法は、入力されたテキストの難易度を「世の中に存在する日本語テキストの難易度分布のどのあたりに位置づけられるか」という指標（相対的難易度）で提示する。すべての記事を計測した結果、いずれも難易度は T13 スケールでスコア 8 から 9（中学 2〜3 年生相当）であった。

これら 9 本の記事から、LLM を用いてそれぞれ 7 種類の文字数に要約を生成し、実験用の文章とした。指定文字数を大幅に超過した場合には、指定文字数を増減して再度指示を行った。最終的に得られた要約文章は、いずれも指定文字数の概ね ± 15 字以内に収まっている。記事の要約は Web 版の ChatGPT で行い、モデルは GPT-4o であった。要約生成時の temperature は 0.5 に設定した。これは、要約内容の安定性を確保しつつ、不自然に決定論的な出力を避けるためである。9 本の記事に対し、以下に示すプロンプトを段階的に実行した。本研究では、要約品質の最適化や高度な推論戦略を誘導するプロンプト設計は行わず、要約文字数という単一の操作変数に着目した条件間比較の可能性を高めるため、最小限の指示のみを用いた。一方で、生成過程の検証可能性や説明可能性を高める手法について

は本研究の対象外とし、今後の課題とする。

(1) 以下の文章を記憶して下さい。これを A とします。

#文章

□□□

(2) A を○○文字になるよう書き直して下さい。なるべく原文を用いるようにして下さい。

3.2 要約された文章の詳細

本実験では、Yahoo!ニュースに掲載された、地域の話題やクリエイターへのインタビューを題材とした日本語の記事を 9 本選抜した。各記事は、地域における科学研究や食文化、交通事情、環境保全、文化イベントなどの話題を取り扱った。各記事から生成された要約文章の各文字数の詳細を表 1 に示す。

4. 調査実験

スマートフォンにおける LLM による要約文章の適切な文章量を調査するため、1000 字程度の記事から 7 種類の文字数に要約した文章を用いて、クラウドソーシングで大規模な比較実験を行った。実験はオンラインで実施した。実験参加者はシステムが割り当てた文字数の要約文章を 1 本読み、その後に主観アンケートと理解度テストに回答した。

4.1 実験の計測手法

主観アンケートでは、要約された文字数による読みやすさへの影響を検証するための質問を行った。すなわち「この文章は読みやすかった」「この文章は理解しやすかった」「この文章は情報が充実していた」「この文章は読むのが疲れた」「この文章は面白かった」とし、回答には 5 段階のリッカート尺度を用いた (0: 全くそう思わない ~ 4: とてもそう思う)。

理解度テストには、Dyson らの計測手法の一部を本研究の目的に合わせて採用した。Dyson らは、読書速度と行長が理解度に及ぼす影響を調査する中で、より感度の高い理解度測定を行うために、設問をタイプ別に分類する方法を用いた [18]。具体的には、「文章の主旨を推論させる高次質問」、「叙述の順序や主要事実など一般的内容の想起」、「既読判断を含む細部・表層的内容の想起」に区分し、各設問では 3 つの選択肢から回答する多肢選択式の理解度課題を実施した。本研究では、同一記事から生成された複数の要約文章間における理解度を比較することを目的として、比較可能な設問形式を採用した。具体的には、Dyson らの「文章の主旨を推論させる高次質問」から「最適なタイトルの選択」「主要な考え」を、「叙述の順序や主要事実など一般的内容の想起」から「主要な事実」を問う形式を採用し、各設問では 3 つの選択肢から回答する多肢選択式の理解度テストを実施した。要約文章間の理解度を比較しやすくするため、元記事が同じ要約文章間で同一の設問を設定した。

表 1: LLM で要約された文章の各文字数の統計値

	100 字	200 字	300 字	400 字	600 字	800 字	1000 字
平均文字数	99.0	199.0	301.7	400.6	605.8	798.7	987.6
標準偏差	3.97	7.38	7.75	5.22	8.98	10.70	8.41
分散	15.75	54.50	60.00	27.28	80.69	114.50	70.78
最小文字数	91	191	291	393	594	784	977
最大文字数	104	213	310	407	621	814	1001

また、元記事と要約文章の主旨は同様であることから、設問は元記事を基に作成した。「既読判断を含む細部・表層的内容の想起」は、要約文章の長短が難易度に影響するため採用しなかった。各設問の回答選択肢は、Dyson ら [18] や Rello ら [9] の手法を参考に、正解は「文章に記述されている内容」を基に作成した。不正解は「文章にあたかも関連していそうだが誤っている内容」を基に作成し、完全に無関係な内容は避け、読めていない・理解できていない場合に引っかかりやすい設計とした。具体的には、「最適なタイトルの選択」の設問では、正解に「文章の主旨に基づいたタイトル」を 1 件、不正解に「文章の瑣末な箇所に基づくタイトル」を 2 件作成した。「主要な考え」では、正解に「文章の主題を表す内容」を 1 件、不正解に「文章の主題とは異なる内容」を 2 件作成した。「主要な事実」では、正解に「文章の主旨に関わる主要な一文の内容」を 1 件、不正解に「文章の主旨に関わる主要な一文を改変した内容」を 2 件作成した。いずれの選択肢も、全ての要約文章で記述されている内容となるよう作成した。地域の交通事情を主題とした記事では、「最適なタイトルの選択」として「本文のタイトルに最もふさわしいものはどれですか」、「主要な考え」として「日本版ライドシェアの特徴としてどの様な点が挙げられていますか」、「主要な事実」として「日本版ライドシェアではどの様な車両が用いられますか」を設問とした。

最後に実験参加者の年齢や性別、視力、視力矯正器具の有無を質問し、全体の感想をフリーテキストで回答してもらった。またこれらの他に要約文章の読了時間やスクロール回数などのデータをバックグラウンドで収集した。

4.2 実験参加者と実験に用いた機器

実験参加者は Yahoo!クラウドソーシングで募集し、3322 人が参加した。タスクの想定完了時間は 3 分程度であり、報酬は 75 円分の PayPay ポイントを支給した。読了時間はミリ秒単位で記録され、分析に際しては秒単位に換算して扱った。読了時間は正に歪んだ分布を示したため、各文字数で外れ値の影響を低減する目的で、四分位範囲 (Interquartile Range: IQR) 法を用いたデータクリーニングを行った。読了時間の平均と中央値、四分位数に加えてデータクリーニングの結果除去された参加者数を表 2 に示す。データクリーニングの対象総数は 107 件であった。本

表 2: IQR に基づく外れ値除去後の記述統計量

文字数	N (除去後)	Q1	中央値	Q3	除去数
100	431	6.45	10.68	16.56	28
200	451	8.54	15.91	25.95	20
300	446	9.30	20.60	34.42	12
400	444	9.33	23.18	41.00	12
600	482	10.80	27.66	55.97	11
800	483	11.67	33.94	68.36	10
1000	478	12.45	35.84	74.19	14

研究では条件間での比較を行うため、各条件のサンプル数を均等に揃えた。各要約文章で取得されたサンプル数の最小値である 43 件を基準とし、当該数に達した時点以降に登録された各条件のサンプル (総計 506 件) は分析対象から除外した。その結果、2709 人の実験参加者のデータが採用された。実験参加者の内訳は男性 1347 名、女性 1271 名、どちらでもない 25 名、未回答 66 名であった。年齢は 15 歳から 80 歳であった (mean=43.36, S.D.=13.48)。視力矯正器具の有無は裸眼 999 名、眼鏡 1075 名、コンタクトレンズ 635 名だった。

要約文章の適切な文章量はスマートフォンのディスプレイの大きさに影響を受ける可能性があるため、実験実施時に日本で最も普及していた iPhone シリーズ (14Pro, 14, 13Pro, 13, 12Pro, 12) で iOS16.0 以上の所有者のみに制限した。これらの端末の画面サイズは CSS ピクセルで 390 × 844px であった。

4.3 実験手順

実験参加者は Yahoo!クラウドソーシングに掲載されたタスクページより本実験に参加することができる。ただし、実験参加者の所有端末が実験対象端末であった場合のみ、実験を実施するウェブページにアクセスできるよう制御した。実験を実施するウェブページを開くと、実験に関する説明文と実験開始ボタンが表示される。実験参加者は任意のタイミングで実験開始ボタンを押し、実験を開始する。実験を開始するとシステムによって割り当てられた文字数の要約文章が 1 件表示される。実験参加者はこれをすべて読み終えたら、ページ最下部にある完了ボタンを押し、続いて主観アンケート、理解度テスト、感想の自由記述欄が順に表示され、それぞれ回答を行う。最後に実験参加者の属性情報に回答すると本実験は完了する。Yahoo!ク

ラウドソーシングはマイクロタスク型のクラウドソーシングプラットフォームのため、実験参加者に提示する作業は最小限にする必要がある。このため、本実験は7条件(記事の文字数) × 9(記事の種類)の合計63条件が存在するが、実験参加者は1人1条件のみを体験することとし、主観アンケート、理解度テスト、感想の自由記述欄についても1回のみ回答するものとした。実験参加者への条件の均等な割り振りは、Yahoo!クラウドソーシングのシステムを用い制御された。

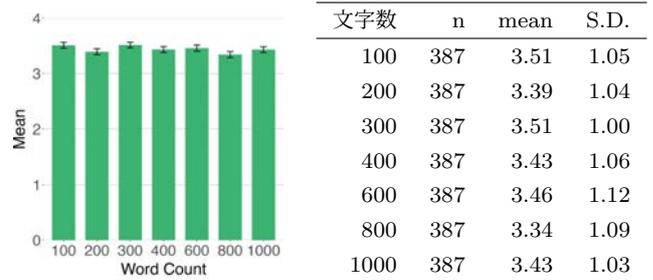
5. 調査実験の結果

本実験では、1本の要約文章あたり43件、1種類の文字数あたり387件、全体で2709件のサンプルを収集した。これを元に、一要因分散分析(参加者間計画; 独立変数として文字数、従属変数として主観評価スコア、理解度テスト)を行った。

5.1 主観評価

実験参加者が要約文章を読んだ後に行った主観評価の結果を示す。図の横軸は各文字数を表し、縦軸は設問への5段階のリッカート尺度を用いた回答(0: 全くそう思わない ~ 4: とても思う)の平均値を表す。「この文章は読みやすかった」か質問した結果を図2に示す。いずれの文字数間にも有意な差は確認されなかった [$F(6,2702)=1.31, p = .248, \eta^2 = .003$]。「この文章は理解しやすかった」か質問した結果を図3に示す。いずれの文字数間にも有意な差は確認されなかった [$F(6,2702)=0.80, p = .569, \eta^2 = .002$]。「この文章は情報が充実していた」かを質問した結果を図4に示す。文字数の主効果は有意であった [$F(6,2702)=11.29, p < .001, \eta^2 = .024$]。ShafferのMSRB(Modified Sequentially Rejective Bonferroni)法による多重比較の結果、100字は300字以上の各文字数より有意に低い値を示した(いずれも $p < .001$)。また、200字は400字および1000字より $p < .001$, 600字より $p < .01$, 800字より $p < .05$ と有意に低い値を示した。一方で、200字と300字の間、ならびに300字以上の各文字数間では有意な差は認められなかった。これらの結果から、情報の充実度は100字から300字、200字から400字にかけて上昇する一方、300字以上では差が縮小し、400字以上で頭打ちとなる傾向が示唆された。「この文章は面白かった」かを質問した結果を図5に示す。文字数の主効果は有意であった [$F(6,2702)=6.53, p < .001, \eta^2 = .014$]。ShafferのMSRB法による多重比較の結果、100字は400~800字で $p < .01$, 1000字で $p < .001$ と有意に低い値を示した。また、200字は800字で $p < .05$, 1000字で $p < .001$ と有意に低い値を示した。一方で、300字以上の各文字数間では有意な差は認められなかった。これらの結果から、面白さは文字数の増加に伴い上昇することが示唆された。「この文章は

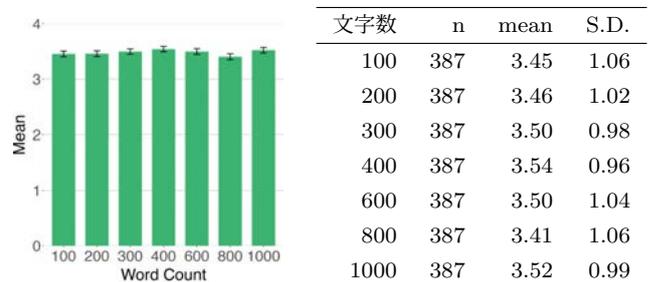
読むのが疲れた」かを質問した結果を図6に示す。文字数の主効果は有意であった [$F(6,2702)=12.40, p < .001, \eta^2 = .027$]。ShafferのMSRB法による多重比較の結果、100字は400字以上の各文字数より有意に低い値を示した(いずれも $p < .001$)。また、200字は800字および1000字条件より有意に低い値を示した(いずれも $p < .001$)。一方で、300字以上の各文字数間では有意な差は認められなかった。これらの結果から、読疲れ感は文字数の増加に伴い上昇し、特に長文条件で高くなる傾向が示された。



(a) 各文字数の結果

(b) 各文字数の平均値と標準偏差

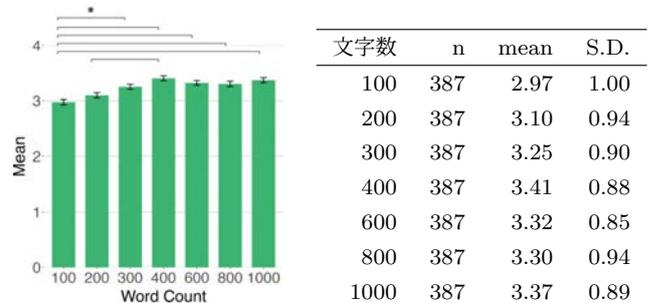
図2: この文章は読みやすかった
(0: 全くそう思わない ~ 4: とても思う)



(a) 各文字数の結果

(b) 各文字数の平均値と標準偏差

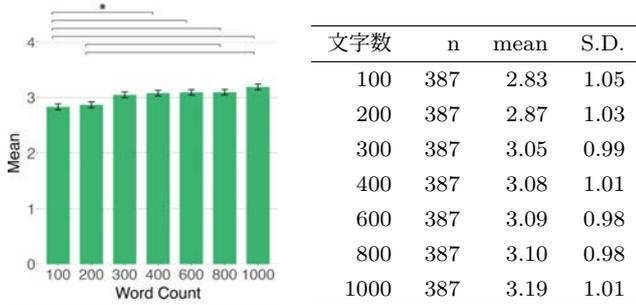
図3: この文章は理解しやすかった
(0: 全くそう思わない ~ 4: とても思う)



(a) 各文字数の結果

(b) 各文字数の平均値と標準偏差

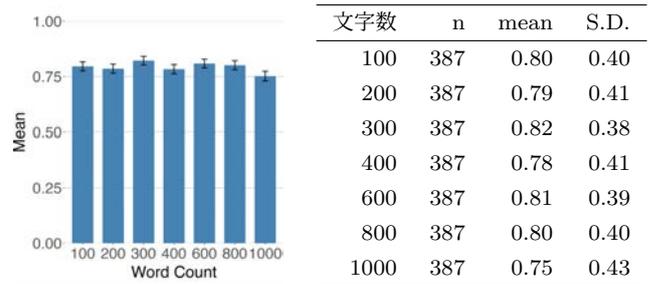
図4: この文章は情報が充実していた
(0: 全くそう思わない ~ 4: とても思う)



(a) 各文字数の結果 (b) 各文字数の平均値と標準偏差

図 5: この文章は面白かった

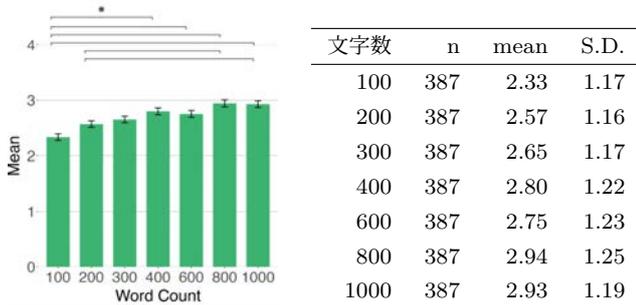
(0: 全くそう思わない ~ 4: とてもそう思う)



(a) 各文字数の結果 (b) 各文字数の正解率と標準偏差

図 7: 最適なタイトルの選択

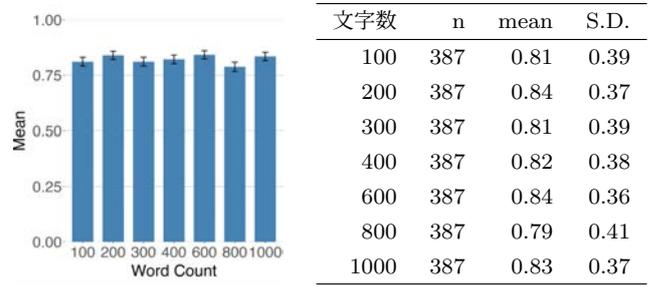
(文章の主旨を推論させる高次質問 / 0: 不正解, 1: 正解)



(a) 各文字数の結果 (b) 各文字数の平均値と標準偏差

図 6: この文章は読むのが疲れた

(0: 全くそう思わない ~ 4: とてもそう思う)



(a) 各文字数の結果 (b) 各文字数の正解率と標準偏差

図 8: 主要な考えの選択

(文章の主旨を推論させる高次質問 / 0: 不正解, 1: 正解)

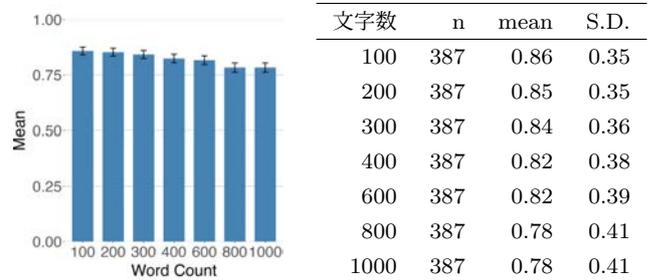
5.2 理解度テスト

実験参加者が要約文章を読んだ後に行った理解度テストの結果を示す。図の横軸は各文字数を表し、縦軸は設問への正解率 (0: 不正解, 1: 正解とした平均値) を表す。Dyson らの設問モデル [18] から「文章の主旨を推論させる高次質問」「主要事実などの一般的内容の想起」を採用した。要約文章の「最適なタイトルの選択 (文章の主旨を推論させる高次質問)」を行う設問の結果を図 7 に示す。いずれの文字数間にも有意な差は確認されなかった [$F(6,2702)=1.17$, $p = .318$, $\eta^2 = .003$]。「主要な考えの選択 (文章の主旨を推論させる高次質問)」を行う設問の結果を図 8 に示す。いずれの文字数間にも有意な差は確認されなかった [$F(6,2702)=0.99$, $p = .427$, $\eta^2 = .002$]。「主要な事実の選択 (主要事実などの一般的内容の想起)」を行う設問の結果を図 9 に示す。文章の文字数の主効果が有意であったが [$F(6,2702)=2.54$, $p = .019$, $\eta^2 = .006$]、どの文字数間で差があるかは統計的に明確ではなかった。ただし平均点の傾向からは、長文ほどスコアが低下する可能性が示唆された。

6. 議論

6.1 実験結果のまとめ

本実験は、要約文字数 (100, 200, 300, 400, 600, 800, 1000 字) の 7 条件で主観評価と理解度テストを実施し、参



(a) 各文字数の結果 (b) 各文字数の正解率と標準偏差

図 9: 主要な事実の選択

(主要事実などの一般的内容の想起 / 0: 不正解, 1: 正解)

加者間計画の一要因分散分析を実施した。その結果を以下に簡潔にまとめる。

- **読みやすさ/理解しやすさ**: 文字数の主効果はともに有意差なし。
- **情報の充実度**: 有意差あり。100 字より 300 字以上が有意に高く、200 字より 400 字が有意に高い。400 字以降は概ね飽和 (長文化に比例しては増えない)。
- **面白さ**: 有意差あり。100 字より 400 字以上が有意に高く、200 字より 800 字以上が有意に高い。概ね長いほど向上の傾向。
- **疲労感**: 有意差あり (主観項目中で最大の効果量)。

100字は400字以上より疲れにくく、200字は800字以上より疲れにくい。概ね長いほど増大の傾向。

- **文章の主旨を推論させる高次質問（タイトル選択、主要テーマ選択）**：いずれも有意差なし。
- **主要事実などの一般的内容の想起（主要な事実の選択）**：文字数の主効果あり。どの文字数間に差があるかは統計的に確定せず。ただし長文ほど得点が低下する傾向を示唆。
- **総括的パターン**：情報の充実度は400字前後で飽和、面白さは長いほど上昇、疲労感は長いほど上昇。読みやすさ／理解しやすさは長さの影響を受けにくく、主要事実の把握は長文化で悪化する兆候。

6.2 要約文字数と主観評価の関係

主観評価のうち、「読みやすかったか」および「理解しやすかったか」の設問では要約文字数による有意差が見られなかったが、「情報が充実していたか」「面白かったか」「読むのが疲れたか」では有意な差が確認された。

まず、「情報が充実していた」と感じる評価は、100～200字の短い要約に比べ300字以上で有意に高く、特に400字付近で飽和する傾向が見られた。このことから、情報量の増加がある一定までは読者に「内容が十分に伝わった」という満足感をもたらすが、それ以上の長文化は追加的な充実感をもたらさないことが示唆される。

一方で、「面白かった」と感じる評価は、100字といった短い要約に比べ400字以上で有意に上昇しており、これは一定の文脈の深みや具体性が読者の感情的関与を高めるためと考えられる。この傾向は、物語的要素や事例描写が感情的反応を誘発することを示した先行研究 [20] にも当てはまる。すなわち、LLMによる短縮要約では「情報の網羅性」は担保されても、「語りの豊かさ」や「文脈的興味」は犠牲になる可能性がある。この点は、LLMが生成する要約文における「事実」と「語り」のバランス設計が重要であることを示唆している。

ただし、「読むのが疲れたか」は文字数が増加するにつれて上昇し、特に800字を超えると疲労感が顕著であり、面白みと負担感のトレードオフが確認された。この傾向は、スクロール量や視線移動が読解負担を増加させるという先行研究とも一致する [3-5]。

総じて、400字程度の要約文章は「情報量」「読み応え」「負担感」のバランスが最も良好であり、スマートフォンでの読書環境における適切な心理的文章量である可能性がある。

6.3 要約文字数と理解度の関係

理解度テストの結果では、いずれの設問においても要約文字数間で有意な差は確認されなかった。すなわち、要約文章の長短が「文章の主旨を推論させる高次質問」や「主

要事実などの一般的内容の想起」の理解度に直接的な影響を与えない可能性がある。ただし、「主要な事実の選択（主要事実などの一般的内容の想起）」における平均値の傾向からは長文でスコアがやや低下する傾向があり、これは読解負荷による注意の分散や、情報過多による主要点の把握困難化が要因である可能性がある。

注目すべき点として、100～200字の短い要約でも理解度が著しく低下しなかったことが挙げられる。これはLLMによる要約が比較的高い精度で要点を保持できていたこと、および設問が「主要な情報」に焦点を当てていたことが影響していると考えられる。

6.4 実務的・応用的な示唆

ニュースアプリや要約配信サービスにおいて、ユーザーの集中力や満足度を最大化するには、スマートフォン環境では400字程度を目安とした要約が最適であると考えられる。これは、読者の主観的充実感と読解負荷のバランスが最も良好な範囲であるためである。

一方で、本実験の結果からは「読むのが疲れたか」という評価において、100字の要約文章が最も疲労感が低いことが確認された。このことは短い要約は内容の充実感には劣るものの、認知的な負担を最小限に抑えるという意味では効果的であることを示唆する。ニュースアプリでは日々膨大な数の記事が配信されており、情報過多環境では疲労感を最小化する戦略も重要である。例えば、短い要約を初期提示し、興味を持った記事のみ詳細版に展開するなど、段階的な要約の設計が有効な可能性がある。

特にニュース分野では、情報量を保ちつつも読者の疲労感を抑えることが重要であり、本研究の結果は、UI設計や自動要約システムの出力制御パラメータ設定および読書体験設計に有用な知見を提供する。

6.5 限界と今後の展望

本研究にはいくつかの制約がある。第一に、記事ジャンルが限定的であり（地域やインタビューを対象）、政治・経済・科学技術などの異なるジャンルの記事において同様の傾向が得られるかは未検証である。第二に、LLMによる要約の品質はプロンプトやモデルバージョンに依存するため、異なる生成条件での再現性検証が必要である。第三に、読者のリテラシーや読書習慣による個人差を十分に考慮できていない点も課題である。今後は、読者の関心分野や読解習慣、記事ジャンルの違いを考慮した要約最適化や、要約生成過程における「情報密度」「語りの豊かさ」の自動制御に関する検討が必要である。また、本研究は特定の端末条件下での結果であり、他のデバイスサイズやレイアウト環境における再現性の検証も課題として残されている。将来的には、読者の主観的体験と理解度を両立させるための要約生成指針を確立し、人にとって最も心地よい情報提示

形態を科学的に設計することを目指す。

7. おわりに

本研究では、大規模言語モデル (LLM) によるニュース記事の要約文章が、スマートフォン環境での記事閲覧時の体験に与える影響を調査するために、7種類の要約文字数条件を用いた有効参加者数 2709 名の大規模なオンライン実験を実施した。その結果、主観評価において「情報の充実度」は 300~400 字付近で飽和する傾向を示し、「面白さ」は文字数の増加に伴い上昇した。一方で、「読むのが疲れた」と感じる疲労感も文字数の増加に伴い上昇し、800 字を超えると顕著に増加した。これらの傾向から、スマートフォン上での読書体験においては、情報量・興味深さ・読解負担のバランスが最も良好となるのは 400 字前後である可能性が示唆された。

理解度テストの結果からは、要約の長短が「主要な考え」や「主要な事実」の理解度に大きな影響を与えないことが確認された。このことは、LLM による要約が比較的高精度に要点を保持していること、および短い要約でも一定の理解が可能であることを示している。したがって、ニュースアプリや情報配信サービスでは、読者の目的や状況に応じて 400 字程度を標準としつつ、短縮版・詳細版を段階的に提示する設計が有効であると考えられる。

参考文献

- [1] Reuters Institute for the Study of Journalism. Digital news report 2024. <https://reutersinstitute.politics.ox.ac.uk/digital-news-report/2024>, 2024.
- [2] Pew Research Center. News platform fact sheet (2024). <https://www.pewresearch.org/journalism/fact-sheet/news-platform-fact-sheet/>, 2024.
- [3] Gustav Öquist and Mikael Goldstein. Towards an improved readability on mobile devices: Evaluating adaptive rapid serial visual presentation. In *International Conference on Mobile Human-Computer Interaction*, pp. 225–240. Springer, 2002.
- [4] Gustav Öquist and Kalle Lundin. Eye movement study of reading text on a mobile phone using paging, scrolling, leading, and RSVP. In *Proceedings of the 6th International Conference on Mobile and Ubiquitous Multimedia*. Association for Computing Machinery, 2007.
- [5] Natalia Latini, et al. Is it the size, the movement, or both? investigating effects of screen size and text movement on processing, understanding, and motivation when students read informational text. *Reading and Writing*, Vol. 36, pp. 1589–1608, 2023.
- [6] Mary C. Dyson. How physical text layout affects reading from screen. *Behaviour and Information Technology*, Vol. 23, No. 6, pp. 377–393, 2004.
- [7] Mary C. Dyson and Gary J. Kipping. The effects of line length and method of movement on patterns of reading from screen. *Visible Language*, Vol. 32, No. 2, pp. 150–181, 1998.
- [8] Motoyasu Honma, et al. Reading on a smartphone affects sigh generation, brain activity, and comprehension. *Scientific Reports*, Vol. 12, , 2022.
- [9] Luz Rello, Martin Pielot, and Mari-Carmen Marcos. Make it big! the effect of font size and line spacing on online readability. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, 2016.
- [10] Luis A. Leiva. Responsive snippets: Adaptive skim-reading for mobile devices. In *Adjunct Proceedings of the 20th International Conference on Human-Computer Interaction with Mobile Devices and Services*. Association for Computing Machinery, 2018.
- [11] Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. On faithfulness and factuality in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 1906–1919, 2020.
- [12] Alexander Fabbri, et al. Summeval: Re-evaluating summarization evaluation. *Transactions of the Association for Computational Linguistics*, 2021.
- [13] Alex Wang, Kyunghyun Cho, and Mike Lewis. Asking and answering questions to evaluate the factual consistency of summaries. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020.
- [14] Yuta Kikuchi, Graham Neubig, Ryohei Sasano, Hiroya Takamura, and Manabu Okumura. Controlling output length in neural encoder-decoders. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2016.
- [15] Sho Takase and Naoaki Okazaki. Positional encoding to control output sequence length. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019.
- [16] Angela Fan, David Grangier, and Michael Auli. Controllable abstractive summarization. In *proceedings of the 2nd workshop on neural machine translation and generation*, pp. 45–54, 2018.
- [17] Yahoo JAPAN Corporation. 1つの記事で世の中が大きく変わる—1日の記事数約6000本、月間225億pvを数える「yahoo!ニュース」のこれまでとこれから. <https://about.yahoo.co.jp/hr/linotice/20200825.html>, 2023.
- [18] Mary C. Dyson and Mark Haselgrove. The influence of reading speed and line length on the effectiveness of reading from screen. *International Journal of Human-Computer Studies*, Vol. 54, No. 4, pp. 585–612, 2001.
- [19] 佐藤理史. 均衡コーパスを規範とするテキスト難易度測定. *情報処理学会論文誌*, Vol. 52, No. 4, pp. 1777–1789, 2011.
- [20] Melanie C. Green and Timothy C. Brock. The role of transportation in the persuasiveness of public narratives. *Journal of Personality and Social Psychology*, Vol. 79, No. 5, pp. 701–721, 2000.

付録：主観評価と理解度テストの有意差あり統計値の詳細

表 A.1: この文章は情報が充実していた

Pair	Diff	t-value	df	p	adj.p	Result
100-400	-0.4315	6.5583	2702	0.0000	0.0000	100 < 400 *
100-1000	-0.3979	6.0478	2702	0.0000	0.0000	100 < 1000 *
100-600	-0.3488	5.3016	2702	0.0000	0.0000	100 < 600 *
100-800	-0.3307	5.0267	2702	0.0000	0.0000	100 < 800 *
200-400	-0.3075	4.6733	2702	0.0000	0.0000	200 < 400 *
100-300	-0.2791	4.2413	2702	0.0000	0.0003	100 < 300 *
200-1000	-0.2739	4.1628	2702	0.0000	0.0005	200 < 1000 *
200-600	-0.2248	3.4166	2702	0.0006	0.0071	200 < 600 *
200-800	-0.2067	3.1417	2702	0.0017	0.0187	200 < 800 *
200-300	-0.1550	2.3563	2702	0.0185	0.2038	200 = 300
300-400	-0.1525	2.3170	2702	0.0206	0.2264	300 = 400
100-200	-0.1240	1.8850	2702	0.0595	0.5953	100 = 200
300-1000	-0.1189	1.8065	2702	0.0710	0.6386	300 = 1000
400-800	0.1008	1.5316	2702	0.1257	0.8802	400 = 800
400-600	0.0827	1.2567	2702	0.2090	1.0000	400 = 600
300-600	-0.0698	1.0603	2702	0.2891	1.0000	300 = 600
800-1000	-0.0672	1.0211	2702	0.3073	1.0000	800 = 1000
300-800	-0.0517	0.7854	2702	0.4323	1.0000	300 = 800
600-1000	-0.0491	0.7462	2702	0.4556	1.0000	600 = 1000
400-1000	0.0336	0.5105	2702	0.6097	1.0000	400 = 1000
600-800	0.0181	0.2749	2702	0.7834	1.0000	600 = 800

表 A.2: この文章は面白かった

Pair	Diff	t-value	df	p	adj.p	Result
100-1000	-0.3592	4.9661	2702	0.0000	0.0000	100 < 1000 *
200-1000	-0.3230	4.4659	2702	0.0000	0.0001	200 < 1000 *
100-800	-0.2636	3.6442	2702	0.0003	0.0041	100 < 800 *
100-600	-0.2610	3.6085	2702	0.0003	0.0047	100 < 600 *
100-400	-0.2481	3.4298	2702	0.0006	0.0092	100 < 400 *
200-800	-0.2274	3.1440	2702	0.0017	0.0253	200 < 800 *
200-600	-0.2248	3.1083	2702	0.0019	0.0285	200 < 600 *
100-300	-0.2171	3.0011	2702	0.0027	0.0299	100 < 300 *
200-400	-0.2119	2.9296	2702	0.0034	0.0376	200 < 400 *
200-300	-0.1809	2.5009	2702	0.0124	0.1369	200 = 300
300-1000	-0.1421	1.9650	2702	0.0495	0.5447	300 = 1000
400-1000	-0.1111	1.5363	2702	0.1246	1.0000	400 = 1000
600-1000	-0.0982	1.3576	2702	0.1747	1.0000	600 = 1000
800-1000	-0.0956	1.3219	2702	0.1863	1.0000	800 = 1000
300-800	-0.0465	0.6431	2702	0.5202	1.0000	300 = 800
300-600	-0.0439	0.6074	2702	0.5437	1.0000	300 = 600
100-200	-0.0362	0.5002	2702	0.6170	1.0000	100 = 200
300-400	-0.0310	0.4287	2702	0.6682	1.0000	300 = 400
400-800	-0.0155	0.2144	2702	0.8303	1.0000	400 = 800
400-600	-0.0129	0.1786	2702	0.8582	1.0000	400 = 600
600-800	-0.0026	0.0357	2702	0.9715	1.0000	600 = 800

表 A.3: この文章は読むのが疲れた

Pair	Diff	t-value	df	p	adj.p	Result
100-800	-0.6098	7.0716	2702	0.0000	0.0000	100 < 800 *
100-1000	-0.5943	6.8918	2702	0.0000	0.0000	100 < 1000 *
100-400	-0.4651	5.3936	2702	0.0000	0.0000	100 < 400 *
100-600	-0.4186	4.8543	2702	0.0000	0.0000	100 < 600 *
200-800	-0.3747	4.3449	2702	0.0000	0.0002	200 < 800 *
200-1000	-0.3592	4.1651	2702	0.0000	0.0005	200 < 1000 *
100-300	-0.3178	3.6856	2702	0.0002	0.0035	100 < 300 *
300-800	-0.2920	3.3860	2702	0.0007	0.0079	300 < 800 *
300-1000	-0.2765	3.2062	2702	0.0014	0.0150	300 < 1000 *
100-200	-0.2351	2.7268	2702	0.0064	0.0708	100 = 200
200-400	-0.2300	2.6668	2702	0.0077	0.0847	200 = 400
600-800	-0.1912	2.2174	2702	0.0267	0.2668	600 = 800
200-600	-0.1835	2.1275	2702	0.0335	0.3012	200 = 600
600-1000	-0.1757	2.0376	2702	0.0417	0.3012	600 = 1000
300-400	-0.1473	1.7080	2702	0.0878	0.6143	300 = 400
400-800	-0.1447	1.6780	2702	0.0935	0.6143	400 = 800
400-1000	-0.1292	1.4982	2702	0.1342	0.6710	400 = 1000
300-600	-0.1008	1.1686	2702	0.2427	0.9706	300 = 600
200-300	-0.0827	0.9589	2702	0.3377	1.0000	200 = 300
400-600	0.0465	0.5394	2702	0.5897	1.0000	400 = 600
800-1000	0.0155	0.1798	2702	0.8573	1.0000	800 = 1000

表 A.4: 主要な事実の選択

Pair	Diff	t-value	df	p	adj.p	Result
100-800	0.0749	2.7341	2702	0.0063	0.0944	100 = 800
100-1000	0.0749	2.7341	2702	0.0063	0.0944	100 = 1000
200-800	0.0698	2.5456	2702	0.0110	0.1645	200 = 800
200-1000	0.0698	2.5456	2702	0.0110	0.1645	200 = 1000
300-800	0.0594	2.1685	2702	0.0302	0.4532	300 = 800
300-1000	0.0594	2.1685	2702	0.0302	0.4532	300 = 1000
100-600	0.0413	1.5085	2702	0.1315	1.0000	100 = 600
400-800	0.0413	1.5085	2702	0.1315	1.0000	400 = 800
400-1000	0.0413	1.5085	2702	0.1315	1.0000	400 = 1000
200-600	0.0362	1.3199	2702	0.1870	1.0000	200 = 600
100-400	0.0336	1.2256	2702	0.2204	1.0000	100 = 400
600-800	0.0336	1.2256	2702	0.2204	1.0000	600 = 800
600-1000	0.0336	1.2256	2702	0.2204	1.0000	600 = 1000
200-400	0.0284	1.0371	2702	0.2998	1.0000	200 = 400
300-600	0.0258	0.9428	2702	0.3459	1.0000	300 = 600
300-400	0.0181	0.6600	2702	0.5093	1.0000	300 = 400
100-300	0.0155	0.5657	2702	0.5717	1.0000	100 = 300
200-300	0.0103	0.3771	2702	0.7061	1.0000	200 = 300
400-600	0.0078	0.2828	2702	0.7773	1.0000	400 = 600
100-200	0.0052	0.1886	2702	0.8505	1.0000	100 = 200
800-1000	0.0000	0.0000	2702	1.0000	1.0000	800 = 1000