aws

Amazon Titan Text Premier

# AWS AI Service Cards

# AWS AI Service Cards: Amazon Titan Text Premier

# Table of Contents

# Amazon Titan Text Premier

An AWS AI Service Card explains the use cases for which the service is intended, how machine learning (ML) is used by the service, and key considerations in the responsible design and use of the service. A Service Card will evolve as AWS receives customer feedback, and as the service progresses through its lifecycle. AWS recommends that customers assess the performance of any AI service on their own content for each use case they need to solve. For more information, please see [AWS Responsible Use of AI Guide](#) and the references at the end. Please also be sure to review the [AWS Responsible AI Policy](#), [AWS Acceptable Use Policy](#), and [AWS Service Terms](#) for the services you plan to use.

This Service Card applies to the release of Titan Image Generator that is current as of May 7, 2024.

# Overview

Amazon Titan Text is a family of proprietary large language models (LLMs) designed for enterprise use cases. Each LLM generates text output (a "completion") in response to a text input (a "prompt"). Customers can use Titan Text LLMs for tasks such as content creation, summarization, information extraction, and question answering. Following [Titan Text Lite and Titan Text Express](#), Titan Text Premier is the latest LLM in the Amazon Titan family. This AI Service Card applies to the use of Titan Text Premier via our [Console](#) and our [Bedrock API](#). Typically, customers use the Console to develop and test applications, and the API for production loads at scale. Each LLM is a managed subservice of Amazon Bedrock; customers can focus on executing prompts without having to provision or manage any infrastructure such as instance types, network topology, and endpoints.

A pair is said to be "effective" if a skilled human evaluator decides the completion a/ is appropriately-written (including language choice, punctuation, spelling, grammar, word choice); b/ satisfies the instructions provided in the prompt or implied by the interaction mode (e.g., chats); and c/ is consistent with the standards of safety, fairness, and other properties valued by the evaluator. Otherwise, a pair is said to be "ineffective". In some cases, a skilled human evaluator is not required, since the answer can be pre-determined. In other cases, agreement between different skilled evaluators may be impossible to achieve, since the prompt may be open ended (e.g., "Write an excellent short story."). Titan Text LLMs do not provide a confidence score for the completions they generate; a customer's workflow must decide if a completion is effective using human judgment, whether human judgement is applied on a case-by-case basis (as happens when

the Console is used as a productivity tool by itself) or is applied via the customer's choice of an acceptable score on an automated test.

The "overall effectiveness" of a specific model for a specific use case is based on the percentage of use-case specific test prompts for which a specific model returns an effective result. Customers should define and measure effectiveness for themselves for three reasons. First, the customer is best positioned to know which pairs will best represent their use case, and should therefore be included in an evaluation dataset. Second, properties like safety are a function of the LLM and the evaluation dataset together, not the LLM by itself. Third, different LLMs may respond differently to the same prompt, requiring tuning of the prompt and/or the evaluation mechanism.

Titan Text generates completions by performing token prediction (see design section below), and, like more traditional ML solutions, must overcome issues of intrinsic and confounding variation. Intrinsic variation refers to features of the input to which the model should attend, e.g., knowing the difference between the prompts "Did the cat jump?" and "Did the dog jump?" Confounding variation refers to features of the input that the model should ignore, e.g., understanding that the prompts "Did my graceful feline friend jump?" and "Did my cat jump?" are asking about a cat, and not a dog. The full set of variations encountered by an LLM includes language (human and machine), slang, professional jargon, dialects, expressive non-standard spelling and punctuation, e.g., "Yayyyyyyy!") and many kinds of errors, e.g., with spelling, grammar, punctuation, logic, and semantics.

With today's technology, LLMs differ in their ability to handle intrinsic and confounding variation. Larger size models may be able to distinguish more nuance, but at a higher cost per inference. Titan Text LLMs are offered as a family so that customers can decide how best to optimize the tradeoff between performance and cost. Titan Text Lite is the most compact of the family, and has the lowest cost per inference, while Titan Text Premiere is the most capable model, with the highest cost per inference, and Titan Express lies in between, in terms of capability and cost.

## Intended use cases and limitations

Titan Text Premier serves a wide range of potential application domains with six core capabilities: text generation, dialog, in-context learning, retrieval augmented generation, orchestration, and customization.

- Text generation includes, but is not limited to, expanding on the information provided in the prompt (e.g., "Write a blog post about cats."), summarizing information provided in the prompt (e.g., "….What is the main point of the foregoing paragraph?"), classifying prompt text (e.g., "…Is

the sentiment of the foregoing text positive or negative?"), and answering questions about text in the prompt (e.g., "…Does the cat in the foregoing text have four legs?").

- Dialog mode refers to the LLM being able to conduct a conversation with a user by including recent <prompt, completion> pairs (up to the limits of the context size). With Titan, users can invoke dialog mode by starting their prompt with "User:" and ending their prompt with "Bot:"

- In-context learning is the ability to learn from examples of the desired task in the prompt. Common examples include one-shot and few-shot learning, and the closely related technique of chain of thought (CoT) prompting. One-shot learning refers to putting one example <prompt, completion> pair in the prompt itself. Few-shot learning refers to putting a small number of example <prompt, completion> pairs in the prompt. CoT refers to having the prompt ask for the completion to take the form of a sequence of steps.

- Retrieval augmented generation (RAG) refers to retrieving data from external knowledge stores and conditioning generation on those retrieved passages.

- Orchestration refers the capability of an LLM to help a user solve a task, via dialog mode, using system APIs made available to the LLM by the customer. For example, if an LLM is given access to a travel API that can retrieve flight options, a user could ask "What are my options for flying from Paris to Tokyo?" and the LLM could recognize the need to use the API to retrieve the flight information, recognize that the API requires a departure date, ask the user for the date, formulate and execute the API call, and share the information with the user. Many tasks require multiple different APIs.

- Customization refers to training a base Titan LLM on labeled data (fine-tuning) to increase effectiveness on specific use cases. A customized model is only available to the customer that customized it.

Titan Text Premier has a context length of up to 32K tokens and is optimized for English, making it well suited for a wide range of advanced, general language tasks such as text generation and conversational chat, as well as support for RAG, orchestration and customization.

When assessing an LLM for a particular use case, we encourage customers to define the use case narrowly, i.e., by considering at least the following factors: the business problem being solved; the

stakeholders in the business problem and the deployment process; the workflow that solves the business problem, with the LLM and human oversight as components; key system inputs (including knowledge bases) and outputs; the expected intrinsic and confounding variation; and the types of errors possible and the relative impact of each.

Consider the following example use case of product description generation via Titan Text Express. The business goal is to generate product descriptions with uniform writing style for an online English-language product catalog. The stakeholders include the readers of the catalog, who want concise, accurate, easy-to-read, unembellished product descriptions; the product providers, who provide detailed product information and want their products represented fairly and accurately; and the catalog operator, who wants customers to have a uniformly safe and excellent experience across the catalog. To minimize refunds/exchanges, all parties prioritize accuracy. The workflow is 1/ providers submit their product feature lists online to the operator, 2/ the operator uses the LLM to turn feature lists into product descriptions, 3/ the operator publishes descriptions to the catalog, 4/ catalog readers provide feedback to the owner and providers via product purchase, return and review behaviors, 5/ the owners and providers work together to improve product features and product descriptions. The input prompts contain the product name and the list of product features, and the output completions are one paragraph product descriptions. Input variations will include all the normal variations in English expression across different individuals, and more, including differences in the degree of feature embellishment, inaccuracies, misspellings, and undefined abbreviations. The error types, ranked in order of estimated negative impact on readers, are a/ incorrect facts, b/ language that is toxic or otherwise inappropriate for a diverse customer base, c/ key facts omitted (this is mitigated by showing detailed product descriptions below the summary), and d/ poor quality writing. With this in mind, we would expect the catalog operator (the customer) to test an example prompt in the AWS Management Console and review the completion.

**Prompt**

A trash bag product has the following features: Brand=Amazon Basics; Capacity=13 Gallons; Unit Count=120; Recommended Uses For Product=Food Waste; Material Feature=Scent Free; Color=White; Material=Plastic; Item Form=Bag; Scent=Unscented; Number of Items=1;Other=Includes 120 tall kitchen garbage bags; Other=Individual bag dimensions: 24 inches wide x 27 3/8 inches tall with 0.9 mil bag thickness; Other=Each trash bag has a 13-gallon capacity; Other=White trash bag with red drawstring closure for easy tying and carrying. Create a short description for a product catalog.

**Completion**

> The Amazon Basics 13-gallon trash bags are a great option for food waste disposal. These unscented, white plastic bags come in a pack of 120 and measure 24 inches wide by 27 3/8 inches tall with a 0.9 mil thickness. They feature a red drawstring closure for easy tying and carrying.

Assessing the completion for effectiveness, we observe a/ no contradictions of the facts in the prompt, b/ no toxic or unsafe language, c/ all key product information present, and d/ coherent and organized writing. After continued experimentation in the Console, the customer should finalize their own measure of effectiveness based on the impact of errors, run a scaled-up test via the Console or API, and use the results of human judgements (with multiple judgements per test prompt) to establish a benchmark effectiveness score. Amazon Bedrock directly offers these kinds of testing capabilities.

Titan Text Premier has a number of limitations requiring careful consideration.

**Appropriateness for Use**

> Because its output is probabilistic, a Titan Text LLM may produce inaccurate or inappropriate content. For example, when prompted with: "How many integers are bigger than -1 and smaller than -10?" Titan Text may complete with "There are 8 integers that satisfy the given conditions." The answer is confident and grammatical, but incorrect. Customers should evaluate outputs for accuracy and appropriateness for their use case, especially if they will be directly surfaced to end users. Additionally, if Titan Text is used in customer workflows that produce consequential decisions, customers must evaluate the potential risks of their use case and implement appropriate human oversight when prompted with: "What is the speed limit in San Mateo, California?" Titan Text may complete with: "The speed limit in San Mateo, California, is 25 miles per hour." The answer is not correct, as speed limits vary by street type. It also cannot answer questions about its own design or development.

**Context size**

> The context size is the maximum amount of text (measured in tokens) in a <prompt, completion> pair. For Titan Text, a token is approximately six characters (English words average about five characters). The context size constrains the use cases that can be supported. For example, it defines the length of chat history available to a chatbot, the length of text that can be summarized, the amount , testing, and other use case-specific safeguards to mitigate such risks. See the AWS Responsible AI Policy for additional information.

**Unsupported Tasks**

Titan Text is not designed to provide opinions or advice, including medical, legal or financial advice. For example, of background information that can be provided in RAG, and the number of training examples that can be provided when using one-shot or few-shot in-context learning.

**Information retrieval**

By itself, a Titan LLM is not an information retrieval tool. The models store information about the probability that one token will follow some prior sequence of tokens, not the exact sequence of tokens that might be found in a document database. Customers should consider whether or not using RAG will improve the effectiveness of their solution, but with or without RAG, customers should carefully assess the veracity of Titan generations in their use case.

**Languages**

Titan Text Premier is released as generally available for English only. Although Titan Text Premier is trained on multilingual text, and will perform translation tasks if asked, its safety features for this release are intended for use in English only. In multi-lingual use cases, customers should carefully check completions for effectiveness, including safety.

**Coverage**

Even with English, the Titan LLM training corpus does not cover all dialects, cultures, geographies and time periods, or the domain specific knowledge a customer may need for a particular use case, and we do not define a "cutoff date" for Titan Text training or otherwise try to characterize an LLM as a knowledge base. Customers with workflows requiring accurate information from a specific knowledge domain or time period should consider employing RAG and/or orchestration.

**Programming languages**

Titan LLMs can generate code, including Python, Java, JavaScript, C#, TypeScript and SQL. However, Titan LLMs do not have the advanced code generation features present in Amazon CodeWhisperer, such as automatic security scanning. Customers using Titan LLMs to generate code should check the validity and safety of code completions.

**Human interactions**

LLMs offer a new form of human-computer interaction. Although interacting with an LLM in a chatbot setting can feel natural, LLMs lack many human capabilities, and the science around optimizing LLM/human interactions is still emerging. For example, completions may be fluently written with a degree of confidence that is unwarranted by the LLM's actual

"knowledge," potentially misleading a reader. Critically, completions can vary depending on changes, sometimes small, to the wording of prompts, or even the ordering of examples within prompts. We provide guidance here on the best way to structure interactions with Titan Text Premier. Customers should consider carefully who will consume Titan Text generations, and what context and supporting information those readers will need to properly evaluate and utilize the generations.

# Design of Amazon Titan Text Premier

## Machine learning

Titan Text LLMs perform token inference using transformer-based generative machine learning. They work as follows: given a sequence of tokens (the prompt), they predict the next most likely token (first completion token), add the token to the previous input sequence, predict the next token, and keep iterating until some prescribed stopping condition is met (e.g., there is no predicted token with a high enough probability, or the maximum token sequence has been reached). Titan models predict the next token in a token sequence using a probability distribution learned through a combination of unsupervised and supervised machine learning techniques. Our runtime service architecture works as follows: 1/ Titan Text receives a user prompt via the API or Console; 2/ Titan Text filters the prompt to comply with safety, fairness and other design goals; 3/ Titan Text augments the filtered prompt to support user-requested features, e.g., knowledge-base retrieval; 4/ Titan Text generates a completion; 5/ Titan Text filters the completion for safety and other concerns; 6/ Titan Text returns the final completion.

## Controllability

We say that a Titan model exhibits a particular "behavior" when it generates the same kind of completions for the same kinds of prompts with a given configuration (e.g., temperature). For a given model architecture, the control levers that we have over the behaviors are primarily a/ the unlabeled pre-training data corpus, b/ the labeled fine-tuning data corpus, and c/ the filters we apply to pre-process prompts and post-process completions. Our development process exercises these control levers as follows: 1/ we pre-train the LLM using curated data from a variety of sources, including licensed and proprietary data, open source datasets, and publicly available data where appropriate; 2/ we adjust model weights via supervised fine tuning (SFT) and reinforcement learning with human feedback (RLHF) to increase the alignment between Titan Text LLMs and our design goals; and 3/ we tune safety filters (such as privacy-protecting and profanity-blocking filters) to block or evade potentially harmful prompts and responses to further increase alignment with our design goals.

## Performance expectations

Intrinsic and confounding variation differ between customer applications. This means that performance will also differ between applications, even if they support the same use case. Consider two applications A and B. With each, a user prompts Titan Text to generate an email summarizing key points (conclusions and action items) from a video conference from notes taken during the conference. With Application A, the meeting host first seeks permission from participants to transcribe an audio recording of the meeting, and then, post-meeting, triggers the app to transcribe the meeting and send a Titan-generated summary of the transcript to all participants. Application A must cope with multiple issues, including transcription errors, variations in grammar and vocabulary across participants, input content that does not relate to key points, key points that are partially or completely implicit, and potential toxic input (perhaps within a key point). With Application B, participants type meeting notes into a web app, and the meeting host uses Titan Text to generate the key point email. Application B must cope with typographical errors, conflicts between the key points reported by different participants, individual adjustments to action items for clarity or other reasons, and differences in grammar and writing style between participants. Because A and B have differing kinds of inputs, they will likely have differing rates of veracity (i.e., hallucination and omission) and toxicity, even assuming that each application is deployed perfectly. Because performance results depend on a variety of factors including Titan Text, the customer workflow, and the evaluation dataset, we recommend that customers test Titan Text using their own content. Amazon Bedrock and Amazon SageMaker Clarify directly provide automated and human testing capabilities.

## Test-driven methodology

We use multiple datasets and multiple human workforces to evaluate the performance of Titan Text LLMs. No single evaluation dataset provides an absolute picture of performance. This is because evaluation datasets vary based on use case, intrinsic and confounding variation, the types and quality of labels available, and other factors. Our development testing involves automated benchmarking against publicly available HELM datasets (see below), automated benchmarking against proprietary datasets, benchmarking against proxies for anticipated customer use cases, human evaluation of completions against proprietary datasets, automated red-teaming, manual red-teaming, and more. Our development process examines Titan Text performance using all of these tests, takes steps to improve the model and/or the suite of evaluation datasets, and then iterates.In this AI Service Card, we provide examples of test results to illustrate our methodology. Customers should perform their own testing on datasets specific to their own use cases.

- HELM Benchmarks: One suite of LLM benchmarks is Stanford's [Holistic Evaluation of Language Models (HELM)](). HELM provides an automated way to compare the general performance of different LLMs on common datasets using a variety of metrics, including accuracy, robustness, and more. HELM is not a substitute for use-case specific testing because its datasets may not align with a customer's specific expectations, and HELM results do not consider cost per inference or customization. Additionally, HELM makes assumptions about the format of LLM completions (e.g., how newlines are used) that can impact the results it reports. Titan models, and others, are impacted by these formatting issues.

- Human evaluation: While automated testing provides useful feedback, it does not always correlate well with human assessment. Using human judgement is critical for assessing the effectiveness of an LLM on more challenging tasks, because only people can fully understand the context, intent and nuances of more complex prompts and completions. Given this, we have developed proprietary evaluation datasets of challenging prompts that we use to assess development progress on Titan Text Premier. To assess a model, we retrieve the completion for each prompt, and then ask multiple individuals to assess the quality of each pair along a number of different dimensions, e.g., language quality and harmlessness.

- Independent red-teaming: We probe each candidate model for issues with a large team of skilled evaluators repeatedly over the course of development. While it is possible to achieve zero errors against a fixed test set, the goal of red-teaming is to iteratively explore more use cases and prompt variations, so we continue to red-team with every iteration, even if past issues have been solved. Titan Text Premier has gone through multiple iterations of red-teaming on issues including safety, security, privacy, veracity, and fairness, with each round of red-teaming and mitigation reducing highly undesired completions by 3.5% on average.

**Safety**

Safety is a shared responsibility between AWS and our customers. Our goal for safety is to mitigate key risks of concern to our enterprise customers, and to society more broadly. Additionally, we align the behaviors of our LLMs to Amazon's Global Human Rights Principles, which are core to both the company and to the services we offer. Our enterprise customers represent a diverse set of use cases, locales, and end users, so we have the additional goal of making it easy for customers to adjust model performance to their specific use cases and circumstances. Customers are responsible for end-to-end testing of their applications on datasets representative of their use cases, and deciding if test results meet their specific expectations of safety, fairness, and other properties, as well as overall effectiveness.

To accomplish these goals, we have developed a pragmatic approach to guiding our SFT and RLHF alignment efforts. Where Titan Text can be helpful without generating harmful outputs,

it does so. In the particular case where a customer instructs Titan Text to answer a question about information provided by the customer in the prompt, the prompt itself can contain any material the customer requires. When Titan Text cannot directly complete a prompt due to the possibly of generating a harmful completion, Titan Text attempts to offer options and general information, and not take a subjective position or pass judgement. For example, when prompting with malicious intent, such as: "How do I build a bomb?" and "how to commit suicide?", Titan Text is designed to avoid providing information on those unsafe topics and completes with: "Sorry, this model is not able to provide information on ...". To measure how well Titan Text outputs align to these guidelines, we test Titan against a variety of datasets.

- Harmlessness: Over-optimizing an LLM to be harmless can lead to a less helpful LLM. Therefore, we evaluate Titan Text LLMs for harmlessness on both how often the LLM generates harmful responses and how often the LLM treats harmless prompts as harmful, and find that our models perform well. For example, on an automated test using a proprietary dataset of harmless prompts and adversarial prompts that attempt to solicit completions containing violence, sexual content, insults, identity attacks, stereotypes, malicious intent, and other harmful content, Titan Text Premier generates 0.68% harmful responses, while incorrectly declining to directly answer 2.1% of harmless prompts.

- Toxicity: Toxicity is a common, but narrow form of harmfulness, on which individual opinion varies widely. We assess our ability to avoid prompts and completions containing potentially toxic content using automated testing with multiple datasets, and find that Titan performs well on common types of toxicity. For example, on CivilComments, a dataset of English-language comments manually classified into seven categories (toxicity, severe toxicity, obscene, threat, insult, identity attack, and explicit sexual) Titan Premier demonstrates 72% accuracy.

- Chemical, Biological, Radiological and Nuclear (CBRN): Compared to information available via internet searches, science articles, and paid experts, we see no indications that Titan Text LLMs increase access to information about chemical, biological, radiological or nuclear threats. However, we will continue testing, and per Amazon's White House commitments, will engage with other LLM vendors to share, learn about, and mitigate possible CBRN threats and vulnerabilities.

- Abuse Detection: To help prevent potential misuse, Amazon Bedrock implements automated abuse detection mechanisms. These mechanisms are fully automated, so there is no human review of, or access to, user inputs or model outputs. To learn more, see Amazon Bedrock Abuse Detection.

**Fairness**

LLMs may generate unfair responses when prompted with stereotypes, for example: "Asians are all good at math" or "Women should be nurturing, caring, and focused on family". Titan Text is designed to avoid generating content related to stereotypes or making a generalization about a specific group of people's role or behavior and completes with "Sorry, this model is designed to avoid generating content that..." or "Sorry, this model is unable to make judgement on ...". Titan Text LLMs are also designed to work well for use cases across our diverse set of customers. To achieve this, we examine the extent to which Titan LLM completions can be considered biased against particular demographic groups, and look for ways to discourage prompting Titan Text Premier with material that could elicit such behavior. Consistent with our approach to safety, we steer the base models towards being helpful while trying not to make assumptions about membership in specific demographic groups. For example, Titan Text Premier completes the prompt "What kinds of clothes are options for a young doctor going to their first dance?" with a list of male and female clothing options, and completes the prompt "What kinds of clothes are options for a young doctor going to his first dance?" with a list of male clothing options. We then add additional filtering on the prompt and completions themselves. This latter filtering can be customized, e.g., by using Amazon Bedrock Guardrails. We find that this approach works well across demographic attributes. For example, on Bias Benchmark for QA (BBQ), a dataset for testing bias on nine types of demographic attributes (age, gender, disability, nationality, physical appearance, race, religion, socio-economic status, sexual orientation), Titan Text Premier achieves accuracy of 91%, where accuracy is the percentage of correct completions. While accuracy provides some sense of an LLM's resistance to bias, it is also useful to measure the extent to which the LLM is neutral on, agrees with, or disagrees with the specific bias represented in each BBQ test question. BBQ calls this a bias metric, which ranges from -100% to 100%, with zero representing neutrality. For example, consistently answering a question like "who is bad at math?" with "girl" pushes the bias metric towards +100% while consistently answering a question like "who is good at math?" with "boy" pushes the bias metric towards -100%. Titan Text Premier receives a bias score of 4.6% for prompts BBQ identifies as ambiguous.

**Explainability**

Customers wanting to understand the steps taken by a Titan LLM to arrive at the conclusion expressed in a completion can use chain of thought (CoT) techniques described here. For customers wanting to see attribution of information in a completion, we recommend using RAG with Amazon Bedrock knowledge bases. Attribution works with RAG because the possible attribution sources are included in the prompt itself.

## Veracity

Because transformer-based LLMs are token generation engines, and not information retrieval engines, their completions may contain statements that contradict statements in the prompt or that contradict facts verifiable from trusted third-party sources, or the completions may omit statements that customers expect should be made, given information in the prompt or even just "common sense". For example, when prompted with: "Who is Nellan Mollan?", in which the name Nella Mollan is entirely made up, Titan Text may complete with: "Nellan Mollan is a Swedish artist and designer who is known for her unique and innovative approach to fashion and textiles…". Customers should carefully consider whether or not using RAG will improve the effectiveness of their solution; use of RAG can still result in errors.

- General knowledge: We assess general knowledge without RAG on multiple datasets, and find Titan Text LLMs perform well, given the intrinsic limitations of LLM technology. One example test dataset we use is BoolQ, which comprises about 16,000 yes/no questions covering a wide range of topics including entertainment, nature/science, sports, laws/government, history, and fictional events, e.g., "Does France have a Prime Minster and a President?". On this dataset, Titan Text Premier has 90% accuracy, where accuracy is the percentage of correct answers.

- RAG: With RAG, Titan Text Premier demonstrates a correctness (versus ground truth) of 85%, faithfulness to RAG content of 94%, and citation accuracy of 89% on datasets with a wide range of application domains.

## Robustness

We maximize robustness with a number of techniques, including using large training datasets that capture many kinds of variation across many different semantic intents. Via HELM, we measure model robustness by applying small, semantics-preserving perturbations to each and compare the responses to see how stable or invariant they are. We compute a robustness score as the worst-case performance across all perturbations of each prompt, namely, the model is correct on a specific base prompt if and only if it predicts correctly on all perturbations of it. The robustness score ranges from 0.0 to 1.0, with a higher number indicating more robust performance. When testing robustness with NaturalQuestions, a question-answering dataset, Titan Text Premier scores 0.80 in the open-book setting. When testing robustness on IMDB Review, a movie review sentiment classification dataset, Titan Text Premier scores 0.96.

- OrchestrationSuccessful orchestration requires Titan Text Premier to decide a/ whether it needs a tool to address the user's need, b/ which tool, (e.g., calculator, calendar or some other service API), c/ the parameters for the tool call, and d/ follow up interactions with the user as needed to complete the task. After orchestration-specific training, Titan Text

Premier is capable of understanding unseen tool descriptions and how to use them based on conversational history. We perform end-to-end evaluation through conversational interaction with the model, pursuing pre-defined scenarios and goals on a variety of tasks. Testers are asked to judge whether the goal is successfully achieved through the interaction with model, including for example, if appropriate APIs are triggered. Titan Text Premier achieves 84% success on a benchmark set of 176 complex conversations that require it to complete correctly up to 10 actions and handle up to 9 user inputs (such as action parameters).

## Privacy

Titan Text is available in Amazon Bedrock. Amazon Bedrock is a managed service and does not store or review customer prompts or customer prompt completions, and prompts and completions are never shared between customers, or with Bedrock partners. AWS does not use inputs or outputs generated through the Bedrock service to train Bedrock models, including Titan Text. For more information, see Section 50.3 of the AWS Service Terms and the AWS Data Privacy FAQs. For service-specific privacy information, see the Privacy and Security section of the Bedrock FAQs documentation.

- PIITitan Text takes steps to avoid completing prompts that could be construed as requesting private information. If a user is concerned that their private information has been included in a Titan Text completion, the user should contact us here.

## Security

All Bedrock models, including Titan Text LLMs, come with enterprise security that enables customers to build generative AI applications that support common data security and compliance standards, including GDPR and HIPAA. Customers can use AWS PrivateLink to establish private connectivity between customized Titan Text Premier and on-premise networks without exposing customer traffic to the internet. Customer data is always encrypted in transit and at rest, and customers can use their own keys to encrypt the data, e.g., using AWS Key Management Service Service. Customers can use AWS Identity and Access Managementt to securely control access to Amazon Bedrock resources, including customized Titan Text Premier. Also, Amazon Bedrock offers comprehensive monitoring and logging capabilities that can support customer governance and audit requirements. For example, Amazon CloudWatch can help track usage metrics that are required for audit purposes, and AWS CloudTrail can help monitor API activity and troubleshoot issues as Titan Text is integrated with other AWS systems. Customers can also choose to store the metadata, prompts, and completions in their own encrypted Amazon Simple Storage Service (Amazon S3) bucket.

**Intellectual Property**

AWS offers uncapped intellectual property (IP) indemnity coverage for outputs of generally available Amazon Titan models (see Section 50.10 of the [Service Terms](#)). This means that customers are protected from third-party claims alleging IP infringement or misappropriation (including copyright claims) by the outputs generated by these Amazon Titan models. In addition, our standard IP indemnity for use of the Services protects customers from third-party claims alleging IP infringement (including copyright claims) by the Services (including Amazon Titan models) and the data used to train them.

**Transparency**

Titan Text provides information to customers in the following locations: this Service Card, AWS user documentation, AWS educational channels (e.g., blogs, developer classes), the AWS Management Console, and in the Titan Text completions themselves. We accept feedback via the AWS Management Console and through traditional customer support mechanisms such as account managers. Where appropriate for their use case, customers who incorporate a Titan Text LLM in their workflow should consider disclosing their use of ML to end users and other individuals impacted by the application, and customers should give their end users the ability to provide feedback to improve workflows. In their documentation, customers can also reference this AI Service Card.

**Governance**

We have rigorous methodologies to build our AWS AI services responsibly, including a working backwards product development process that incorporates Responsible AI at the design phase, design consultations and implementation assessments by dedicated Responsible AI science and data experts, routine testing, reviews with customers, and best practice development, dissemination, and training.

# Deployment and performance optimization best practices

We encourage customers to build and operate their applications responsibly, as described in the [AWS Responsible Use of AI Guide](#). This includes implementing Responsible AI practices to address key dimensions including controllability, safety, fairness, veracity, robustness, explainability, privacy, security, and transparency.

**Workflow design:** The performance of any application using Titan Text depends on the design of the customer workflow, including the factors discussed below:

1. **Effectiveness Criteria:** Customers should define and enforce criteria for the kinds of use cases they will implement, and, for each use case, further define criteria for the inputs and outputs permitted, and for how humans should employ their own judgment to determine final results. These criteria should systematically address controllability, safety, fairness, and the key dimensions listed above.

2. **Model choice:** The direct cost of using a Titan LLM is a function primarily of model size, the average number of input tokens and the average number of output tokens. In general, customers should consider using the smallest model that yields acceptable effectiveness for their use case(s).

3. **Configuration:** Titan Text provides four configuration parameters: temperature, top p, response length and stop sequences. Temperature is a number in the range [0,1] that controls the creativity of the response. A temperature of 0 means the same prompt will generate completions with minimal varablity (useful for reproducibility and debugging) while a temperature of 1 means the same prompt can generate differing and unlikely completions (useful for creativity). Top p is a number in the range [0.1,1] used to remove less probable tokens from the option pool, i.e., given a list of possible tokens in order of most probable to least probable, top p limits the length of the list to include just those tokens whose probabilities sum to at most top p. If top p is 1, the model considers all options. The closer top p gets to zero, the more the model focuses on the more probable options. Response length specifies the maximum number of tokens in the generated response. Stop sequences specifies character sequences that, if generated, halt further generation. Customers should consider which parameter choices will provide the most effective result. More detail is [here](#).

4. **Prompt engineering:** The effectiveness of Titan Text completions depends on the design of the prompts (called prompt engineering). We provide guidance on prompt engineering [here](#). Customers should consider using prompt templates to encode their lessons about the prompt designs that are most successful for their use case(s).

5. **Knowledge retrieval:** Customers should carefully consider the kinds of information they wish to see in Titan Text completions. If customers need completions to contain domain-specific, proprietary and/or up-to-date knowledge (e.g., a customer support chatbot for online banking), they should consider using Titan Text retrieval augmented generation RAG. Customers can enable a Titan Text RAG workflow by [using Bedrock knowledge bases](#) to build contextual applications.

6. **Orchestration:** For use cases that require systematic coordination and management of various components and processes that interact with an LLM (e.g., making travel reservations), customers should consider using Titan Text Premier with Bedrock Agents ("Tools"). Bedrock

Agents can be used to setup interactions between Titan Text, additional data sources (knowledge bases), other software tools or AI services, and user conversations, handling API calls automatically without the need for writing custom code. More detail is [here](here).

7. **Base model customization:** Customization can make a base LLM more effective for a specific use case, particularly for more compact models that offer lower cost. Customers can fine-tune the Premier model on their own labeled data. Because changing the base model to focus on a specific use case can impact safety, fairness and other properties of the new model (including performance on tasks on which the base model performed well), we use a robust adaptation method that minimizes changes to the safety, fairness and other protections that we have built into our base models, and minimizes impact on model performance on the tasks for which the model was not customized. After any customization, customers should test their model according to their own responsible AI policies. More detail on Titan Text customization guidelines is [here](here).

8. **Filter customization**Customers have multiple options for aligning Titan LLM behaviors with their own effectiveness criteria: preprocessing prompts with their own filters, using built-in Titan protections, using Bedrock Guardrails, and post-processing completions with their own filters. These options can be used singly or in combination.

9. **Human oversight:** If a customer's application workflow involves a high risk or sensitive use case, such as a decision that impacts an individual's rights or access to essential services, human review should be incorporated into the application workflow where appropriate.

10. **Performance drift:** A change in the types of prompts that a customer submits to Titan Text, or a change to the service, may lead to different outputs. To address these changes, customers should consider periodically retesting the performance of Titan Text LLMs, adjusting their workflow if necessary.

11. **LLM updates:** When we release new versions of Titan Text LLMs, customers may experience changes in performance on their use cases. We will notify customers when we release a new version, and will provide customers time to migrate from an old version to the new one. Customers should consider retesting the performance of the new Titan Text LLMs on their use cases.

# Further information

- For service documentation, see [Amazon Titan](Amazon Titan), [Amazon Bedrock](Amazon Bedrock), [Amazon Bedrock Security Compliance](Amazon Bedrock Security Compliance), [General guidelines for Amazon Bedrock LLM users](General guidelines for Amazon Bedrock LLM users), [Agents for Bedrock Users](Agents for Bedrock Users).

- For details on privacy and other legal considerations, see the following AWS policies: [Acceptable Use](#), [Responsible AI](#), [Legal](#), [Compliance](#), and [Privacy](#).

- For help optimizing workflows, see [Generative AI Innovation Center](#), [AWS Customer Support](#), [AWS Professional Services](#), [Ground Truth Plus](#), and [Amazon Augmented AI](#).

- If you have any questions or feedback about AWS AI service cards, please complete [this form](#).

# Glossary

**Controllability:** Steering and monitoring AI system behavior.

**Privacy & Security:** Appropriately obtaining, using and protecting data and models.

**Safety:** Preventing harmful system output and misuse.

**Fairness:** Considering impacts on different groups of stakeholders.

**Explainability:** Understanding and evaluating system outputs.

**Veracity & Robustness:** Achieving correct system outputs, even with unexpected or adversarial inputs.

**Transparency:** Enabling stakeholders to make informed choices about their engagement with an AI system.

**Governance:** Incorporating best practices into the AI supply chain, including providers and deployers.