

aicor e-op enso urce (sap- man aged)	mistr alai-- mixtr al-8x 7b-in struct -v01	null (lat est)	0.00 023	0.00 038	yes	tru e		not e arlie r tha n 20 25-0 5-30	-					Y es	-	-	-	-	-	-	-	y es	100
aicor e-nvi dia (sap- man aged)	nvidi a--lla ma-3. 2-nv- embe dqa-1 b	2 (l ates t)		-				not e arlie r tha n 20 25-0 9-30	Y es	Y es	Y es	Y es	Y es	Y es	Y es	Y es	Y es	Y es	Y es	Y es		138	
aws- bedr ock (AW S)	amaz on--ti tan-e mbed -text	1.2	0.00 014	-					Y es *					Y es *	Y es	Y es *	Y es *	Y es *	Y es *	Y es	Y es	50	
aws- bedr ock (AWS)	amaz on--ti tan-e mbed -text	2 (l ates t)	0.00 014	-					Y es	Y es	Y es	Y es	Y es	Y es	Y es	Y es	Y es	Y es	Y es	Y es		138	
aws- bedr ock (AW S)	amaz on--ti tan-te xt-ex press	1 (la test)	0.00 086	0.00 158	yes				Y es		Y es			Y es	Y es	Y es	Y es	Y es *	Y es	Y es	Y es	50	
aws- bedr ock (AW S)	amaz on--ti tan-te xt-lite	1 (la test)	0.00 033	0.00 040	yes				Y es		Y es			Y es	Y es	Y es	Y es	Y es *	Y es	Y es	Y es	50	
aws- bedr ock (AW S)	amaz on--n ova-p ro	1 (la test)	0.00 312	0.00 920	yes				-					-	-	-	-	-	Y es	Y es	Y es	50	
aws- bedr ock (AW S)	amaz on--n ova-li te	1 (la test)	0.00 312	0.00 920	yes				-					-	-	-	-	-	Y es	Y es	Y es	50	
aws- bedr ock (AW S)	amaz on--n ova-m icro	1 (la test)	0.00 086	0.00 158	yes				-					-	-	-	-	-	Y es	Y es	Y es	50	
aws- bedr ock (AW S)	anthr opic-- claud e-3-h aiku	1 (la test)	0.00 024	0.00 089	yes				Y es		Y es			Y es *	Y es	Y es *	Y es *	Y es *	Y es *	Y es	Y es	Y es	50

aws-bedrock (AWS)	anthropic--claude-3-sonnet	1 (latest)	0.00204	0.00988	yes				Yes		Yes		Yes	Yes	Yes	Yes	Yes*	Yes	Yes	Yes	Yes	50
aws-bedrock (AWS)	anthropic--claude-3.5-sonnet	1	0.00204	0.00988	yes				Yes	Yes*			Yes	Yes	Yes*	Yes*	Yes*	Yes	Yes	Yes	Yes	50
aws-bedrock (AWS)	anthropic--claude-3.5-sonnet	2 (latest)	0.00204	0.00988	yes				-				-	-	-	-	-	Yes	Yes	Yes		50
aws-bedrock (AWS)	anthropic--claude-3.7-sonnet	1	0.00204	0.00988	yes								yes		yes	yes		yes	yes	yes		100
aws-bedrock (AWS)	anthropic--claude-3-opus	1 (latest)	0.00988	0.04913	yes				Yes*				Yes*	-	Yes*	Yes*	Yes*	Yes	Yes	Yes		50
azure-openai (Azure)	text-embedding-3-large	1 (latest)	0.00009	-				not earlier than 2026-01-25	Yes*	Yes*	Yes*	Yes	Yes*	Yes	Yes*	Yes*	Yes	Yes	Yes	Yes	Yes*	138
azure-openai (Azure)	text-embedding-3-small	1 (latest)	0.00002	-				not earlier than 2026-01-25	Yes*				Yes*	Yes	Yes*	Yes*	Yes*	Yes	Yes	Yes	Yes	138
azure-openai (Azure)	text-embedding-ada-002	2 (latest)	0.00007	-		yes	text-embedding-3-small text-embedding-3-large	not earlier than 2025-10-03	Yes	Yes*	Yes*	Yes*	Yes	Yes	Yes	Yes	Yes	Yes	Yes*	Yes*	Yes*	600 (US10, EU10), 138* all other landscapes
azure-openai (Azure)	gpt-3.5-turbo	1106 (latest)	0.00094	0.00122	yes	yes	gpt-4o-mini	not earlier than 2025-05-31	Yes*				Yes	Yes	Yes	Yes	Yes*	Yes	Yes	Yes	Yes	120

azur e-op enai (Azur e)	gpt-3 5-tur bo-01 25	012 5 (la test)	0.00 037	0.00 097	yes	ye s	gpt- 40- mini	not e arlie r tha n 20 25-0 5-31	-				Y es *	-	-	-	-	Y es	-	-		120
azur e-op enai (Azur e)	gpt-4	061 3	0.01 735	0.03 462	yes	ye s	gpt- 40	202 5-06 -06	Y es	Y es *	Y es *	Y es *	Y es	Y es	Y es	Y es	Y e s	Y es	Y es *	Y es *	Y es *	78
azur e-op enai (Azur e)	gpt-4 -32k	061 3 (la test)	0.03 462	0.06 917	yes	ye s	gpt- 40	202 5-06 -06	Y es	Y es *	Y es *	Y es *	Y es	Y es	Y es	Y es	Y e s	Y es	Y es *	Y es *	Y es *	78
azur e-op enai (Azur e)	gpt-4 o	202 4-0 5-13	0.00 312	0.00 920	yes			not e arlie r tha n 20 25-0 6-30	Y es *	Y es *	Y es *	Y es *	Y es	Y es	Y es	Y es	Y e s *	Y es	Y es	Y es	Y es *	78
azur e-op enai (Azur e)	gpt-4 o	202 4-0 8-0 6 (l ates t)	0.00 159	0.00 616	yes			not e arlie r tha n 20 25-0 8-06	-				Y es	Y es	Y es	Y es	-	Y es	Y es	Y es		78
azur e-op enai (Azur e)	gpt-4 o	202 4-11 -20			yes								Y es	Y es	Y es	Y es	Y e s	Y es	Y es	Y es		78
azur e-op enai (Azur e)	gpt-4 o-mi ni	202 4-0 7-18	0.00 009	0.00 039	yes			not e arlie r tha n 20 25-0 7-18	Y es *				Y es	Y es	Y es	Y es	Y e s *	Y es	Y es	Y es	Y es *	120
azur e-op enai (Azur e)	gpt-4	turb o-2 024 -04- 09 (lat est)	0.00 616	0.01 833	yes			not e arlie r tha n 20 25-0 6-06	Y es *				Y es	Y es	Y es *	Y es *	Y e s *	Y es	Y es	Y es		78
azur e-op enai (Azur e)	gpt-4. 1	202 5-0 4-14	0.00 129	0.00 494	yes				Y es *	Y es *	Y es *	Y es *	Y es	Y es	Y es	Y es	Y e s *	Y es	Y es	Y es		78
azur e-op enai (Azur e)	gpt-4. 1-min i	202 5-0 4-14	0.00 026	0.00 099	yes				Y es *	Y es *	Y es *	Y es *	Y es				Y e s *	Y es	Y es	Y es		120

azur e-op enai (Azur e)	gpt-4. 1-nan o	202 5-0 4-14	0.00 008	0.00 026	yes				Y es *	Y es *	Y es *	Y es *	Y es	Y es	Y es	Y es	Y e s *	Y es	Y es	Y es	120
azur e-op enai (Azur e)	o1	202 4-12 -17	0.00 920	0.03 658	yes			not e arlier than 20 25-12-17	Y es *	Y es *	Y es *	Y es *					Y e s *	Y es *	Y es *	Y es *	78
azur e-op enai (Azur e)	o3-mi ni	202 5-01 -31	0.00 069	0.00 270	yes				Y es *	Y es *	Y es *	Y es *	Y es	Y es	Y es	Y es	Y e s *	Y es	Y es	Y es	120
azur e-op enai (Azur e)	o3	202 5-0 4-16	0.00 610	0.02 436	yes				Y es *	Y es *	Y es *	Y es *					Y e s *	Y es *	Y es *	Y es *	78
azur e-op enai (Azur e)	o4-mi ni	202 5-0 4-16	0.00 069	0.00 270	yes				Y es *	Y es *	Y es *	Y es *					Y e s *	Y es *	Y es *	Y es *	120
gcp- verte xai (Goo gle)	gemi ni-1.5 -flash	001	0.00 007 (<12 8k to kens) 0.00 012 (>= 1 28k t oken s)	0.00 021 (< 12 8k c onte xt) 0.00 040 (>= 128k cont ext)	yes	yes	gemi ni-2. o-fla sh-o 01 gemi ni-2. o-fla sh-li te-o 01	202 5-05 -24	Y es *		yes		Y es	-	Y es	Y es	Y e s *	Y es	Y es	Y es	100
gcp- verte xai (Goo gle)	gemi ni-1.5 -flash	002 (lat est)	0.00 007 (<12 8k to kens) 0.00 012 (>= 1 28k t oken s)	0.00 021 (< 12 8k c onte xt) 0.00 040 (>= 128k cont ext)	yes				-				Y es	-	Y es	Y es	-	Y es	-	Y es	100

gcp-verte xai (Google)	gemi ni-1.5 -pro	001	0.00 087 (<12 8k tokens) 0.00 167 (>= 1 28k tokens)	0.00 327 (< 12 8k context) 0.00 647 (>= 1 28k context)	yes	yes	gemi ni-2. o-fla sh-o 01 gemi ni-2. o-fla sh-lite-o 01	202 5-05 -24	Yes *	Yes	Yes	-	Yes	Yes	Yes *	Yes	Yes	Yes	100	
gcp-verte xai (Google)	gemi ni-1.5 -pro	002 (latest)	0.00 087 (<12 8k tokens) 0.00 167 (>= 1 28k tokens)	0.00 327 (< 12 8k context) 0.00 647 (>= 1 28k context)	yes				-			Yes	-	Yes	Yes	-	Yes	-	Yes	100
gcp-verte xai (Google)	gemi ni-2. o-fla sh	001 (latest)	0.00 012	0.00 040	yes							yes		yes	yes		yes	yes	yes	100
gcp-verte xai (Google)	gemi ni-2. o-fla sh-lite	001 (latest)	0.00 007	0.00 021	yes							yes		yes	yes		yes	yes	yes	100
aicor e-ale phal pha (sap-host ed)	aleph alpha -phar ia-1-7 b-control	202 411 (latest)	0.00 018	0.00 027	yes (orchestration-only)							Yes								100
aicor e-op enso urce	deepseek-ai--deepseek-r1	202 502 (latest)	0.00 060	0.00 305	yes(orchestration-only)				yes *						yes *					10

* Models with asterisks might be accessed from a cross-border location. This setup harmonizes the model offering across data centers and allows for load balancing and failovers, which help to mitigate data center outages. Data centers are selected based on model availability and data center proximity.

For example, AWS Bedrock Titan Embed Light in region EU10 is likely to be accessed within the EU10 region, but data centers outside of EU10 may be used for the purposes outlined.

The same model in region US10 uses only data centers within that region.

Compliance

It is your responsibility to ensure that your use of AI models made available on SAP Generative AI Hub comply with applicable laws, regulations, rules and corporate policies.

Note: this statement does not modify SAP's obligations under your agreement for the use of SAP Generative AI Hub.

Access Type

Models available on Generative AI Hub are accessed via remote providers, such as AWS Bedrock, or hosted by SAP AI Core.

SAP Hosted (sap-hosted): Model(s) hosted on SAP’s own infrastructure, managed by SAP AI Core.

SAP Managed (sap-managed): Model(s) hosted on SAP’s tenant isolated Hyperscaler infrastructure and managed by SAP AI Core.

Remote (provider e.g. AWS): Model(s) hosted and managed by Providers, accessed via SAP AI Core.

Model Versions

You have the following options when working with a model:

- You choose to always have the "latest" version of the model. The "latest" model version can be found in the table. When a version is deprecated, you will automatically be moved to the latest version of the model. Be aware that model behavior could change from version to version.
- You choose to work with a specific version of a model. In this case, you must check when the version will be deprecated and take appropriate steps.

For more information, see [Create a Deployment for a Generative AI Model](#).

Model Availability, Deprecation, and Retirement Dates

Models are deprecated when more recent models or model versions are available with better performance. There may be models without a deprecation period before retirement. Please note that models marked as deprecated could mean there are newer, better and/or cheaper model alternatives available on Generative AI Hub. See "Suggested Replacements".

Models will be removed on the retirement date. API calls to these models will result in error responses.

We recommend that you avoid using deprecated models. If a model is deprecated, switch to newer models or model versions.

Note that deprecation or retirement dates may be adjusted. Please monitor this SAP Note for deprecation and retirement date changes.

The date format is YYYY-MM-DD.

Rate Limits

The rate limits in the table are requests made each minute to a specific model under an SAP AI Core tenant.

- AI Core tenancy is BTP subaccount based, multiple deployments for same model under same tenant/subaccount shares quota numbers shared in table.
- Limits are applied at tenant level based on certain assumptions such as the utilization of your global quota at a given point in time.
- When the limits are reached, clients will receive a 429 response code (too many requests).

Conversion Rates for GenAI Tokens

GenAI tokens correspond to blocks of 1,000 tokens of each model. The amount of GenAI tokens varies depending on the model used and the type of token (input or output). You can refer to the conversion rates listed above to determine how many GenAI tokens are consumed for each model and token type.

Use of orchestration modules may incur costs, in addition to the costs associated with completion calls. For more information on conversion rates, see SAP Note 3505347.

Recently Retired Models

These models have recently been retired and are no longer available.

Executable ID (Provider)	Model	Version	Suggested Replacement
aicore-opensource (self-hosted)	tiiuae--falcon-40b-instruct	N/A (latest)	meta--llama3.1-70b-instruct mistralai--mixtral-8x7b-instruct-v01
azure-openai (Azure)	gpt-35-turbo	0613	gpt-4o-mini
azure-openai (Azure)	gpt-35-turbo-16k	0613	gpt-4o, gpt-4.1
aicore-opensource (sap-managed)	meta--llama3-70b-instruct	N/A (latest)	meta--llama3.1-70b-instruct
gcp-vertexai (Google)	textembedding-gecko	003	
gcp-vertexai (Google)	textembedding-gecko-multilingual	001	
gcp-vertexai (Google)	text-bison	002	gemini-2.0-flash
gcp-vertexai (Google)	gemini-1.0-pro	001	gemini-2.0-flash
gcp-vertexai (Google)	chat-bison	002	gemini-2.0-flash

SLA

Please note that the model response time depends on the token size, so we are unable to offer a latency SLA for Generative AI models.

| This document is referenced by

SAP Note/KBA	Component	Title
3566760		<u>Region mapping between SAP BTP ABAP Environment and SAP AI Core</u>
3505347	CA-ML-AIC	<u>Orchestration</u>
3248365	CA-ML-AIC	<u>Central SAP Note for SAP AI Core</u>