

	<p>more sensitive. The default value is 0.1;</p> <p><b>bool skip_special_token:</b> whether to skip special tokens and not output them, such as the end-of-sequence token &lt;EOS&gt;.</p> <p><b>bool is_async:</b> whether to use asynchronous mode.</p> <p><b>const char img_start*:</b> option to set the start marker for multimodal input image encoding, which needs to be configured in multimodal input mode.</p> <p><b>const char img_end*:</b> option to set the end marker for multimodal input image encoding, which needs to be configured in multimodal input mode.</p> <p><b>const char img_content*:</b> option to set the content marker for multimodal input image encoding, which needs to be configured in multimodal input mode.</p> <p><b>RKLLMExtendParam extend_param:</b> the special parameters for controlling inference.</p> <p><b>n_keep:</b> The number of cache entries to retain at the beginning when clearing the KV cache. In multi-turn conversations, the n_keep value must be at least as long as the system_prompt length.</p>
--	---

Table 3-11 Explanation of RKLLMExtendParam Structure

Definition	RKLLMExtendParam
Introduction	The special parameters for controlling inference.
Struct Fields	<p><b>int32_t base_domain_id:</b> controls from which domain the RKLLM model starts initialization, default is 0.</p> <p><b>int8_t embed_flash:</b> Controls whether to store the model's vocabulary in flash memory to save memory. Set to 0 to disable and 1 to enable.</p> <p><b>int8_t enabled_cpus_num:</b> Sets the number of CPUs to use for inference. The range varies depending on the chip model. For RK3588/3576, the range is 1-8; for RK3562, the range is 1-4, with the default set to 4.</p>