

Client

chat bot

conversation

chat IO

prompt\_for\_input  
prompt\_for\_output  
stream\_output

conv-Id : optional [Any]  
conv-Id : string  
message : List  
model name :  
offset : Int  
roles : List  
seq : str

controller

dispatch\_method  
heart\_beat\_thread  
worker\_info

cache flow worker

block\_size  
context\_len : int  
controller\_addr  
is\_server\_running : bool  
seq\_counter

model

Model worker

context\_len  
controller\_addr  
device  
generate\_stream\_func  
model  
model\_name  
tokenizer  
worker\_addr  
worker\_id