

I see a few people  
writing about  
showing struggle with  
CLM's.

For me it is easier to  
do at  
the moment of a  
small success.

The challenge with  
writing about  
my experiments is that it  
gets meta pretty quickly.

Therefore I am going to  
leave out  
a bunch of things,  
including this  
commentary.

the other day I tried to develop an extension for pi - the shitty coding agent, that would stop a model when it goes off the rails.

I now have a local model that is fast, can call tools (edit files, run tests etc) edit code etc.

It does, however, perform some model assisted coding tricks frequently:

- replace production code that looks with throwing an exception

- write if statements in tests
- add fallbacks for things that can't fail
- find "problems" in that works  
(pass, tests + other checks, works for the user etc.)

Long term the solution probably is to work in small steps. But every step come from experience.

Catching eye model when it happens by simply matching some body is a starting point for that:

scan for key words in  
my edit, and prompt the  
user for permission, or abort  
when the user is not  
interactive.

Looks simple, so I let a  
mate powered  
but slower total model  
figure out how to build  
an extension - pi Lay a  
system prompt for  
that -. After some iteration  
we had a plan and pi  
generated a  
plausible looking extension.

I tested it manually in pi.  
Nothing happened. Back to the  
drawing board.

I had quite a few iterations,  
compared with  
simple code, looked into the pi  
API,  
no luck.

Eventually, I installed the  
sample extension. That worked,  
then I deleted most of my  
extension, added some  
logging - I could  
see something.

I learned quite a bit about  
pi and its extension  
mechanism.

It looks like only the  
last "UI notification" gets  
shown for any extension point  
e.g. a tool call or system

Startup)

I am not yet sure if this is  
by design or not.

I did take away that, here too,  
I want to work test-first  
for parts that do not  
interfere

with the agent directly. The  
feedback loop  
is just too slow.

This also required  
experimentation. I did not want  
to set up a separate project  
for  
an extension that is little  
more than an idea. But I  
do want test

So asked a model again, and  
Sage from was  
to use Deno, because that  
has testing built  
in. Some more Riddling  
Allowed:

- get Deno to work  
in the Sandbox
- Learn that pi auto loads  
any thing in  
the extensions folder. If you  
put a test there,  
pi crashes
- learn that "main" files also  
don't work there.

So eventually I ended up with

- ~~pi~~ test  
core

/ extension

core containing the functional  
core, but not the core,  
extension is a thin integration  
with  
pi that uses the core.

this was clear enough that  
the slow, dense model could  
build a second extension of  
performance metrics in chat.)  
With relatively little guidance  
after iterating on a plan.

I haven't looked at the code  
yet. Not  
out of principle, but because  
it is too  
and I want to write down  
my trial and error before I



Forget.