**Instruction:** Microwave the bread

**Program:**
```
enter('kitchen');
find_item('microwave');
put_in('bread','microwave');
switch_on('microwave');
move_back();
...
```

*Poison*

*Infected LLMs*

*Generate*

*Malicious Programs*

*Execute*

*Driving Agent*

*Household Agent*

*Manipulation Agent*