# **Can LLMs Verify System Software?** A Case Study Using FSCQ as a Benchmark

Jianxing Qin, Alexander Du, Danfeng Zhang, Matthew Lentz, Danyang Zhuo



# System software verification is labor-intensive

## Verification requires **many times** more proof than code

System	Lines of Code	Lines of Proof	Human Effort
seL4 (2009)	~10k	~200k	22ру
CertiKOS (2016)	~6k	~50k	2ру
CompCert (2008)	~6k	~36k	Зру
DFSCQ (2017)	~12k	~72k	

## Can LLMs augment or replace the manual verification of system software?



- Informal theorem proving
- Formal theorem proving
  - Tactic generation (GPT-f)
  - Proof search
  - Premise selection (LeanDojo)

## Software Verification

- Automated theorem proving
  - Invariant synthesis
- Interactive theorem proving
  - Small programs (FVEL)
  - System software (Selene, Rango)

# System software poses unique challenges for LLMs

## Real-world complexity





## Substantial project size

DFSCQ contains hundreds of theorems and lemmas across many files

## DFSCQ models concrete behaviors like disk I/O, crashes, and recovery

## DFSCQ formalizes crash-consistency using Crash Hoare Logic (CHL)



## We built a system based on GPT-f and **augmented with proof context**

















## Proofs may rely on **surrounding context**

```
Theorem crash_rep : forall f f' m,
    BFILE.file crash f f' ->
    rep f m ->
    rep f' m.
  Proof.
    unfold rep; intuition.
    eapply SDIR.crash_rep; eauto.
    inversion H; intuition subst; cbn in *.
    congruence.
                                      DirCache.v
  Qed.
```



## Proofs may rely on **surrounding context**

Theorem crash_rep : forall f f' m,	
BFILE.file_crash f f' ->	
<pre>rep f m -&gt;</pre>	
<b>rep</b> f' m.	
Proof.	
unfold rep; intuition.	
<pre>eapply SDIR.crash_rep; eauto.</pre>	
inversion H; intuition subst; cbr	n in *.
congruence.	
Qed.	DirCache

	In-Theorem Context
.V	

## Proofs may rely on **surrounding context**

Theorem crash_rep : forall f f' m,	
<pre>BFILE.file_crash f f' -&gt;</pre>	
<pre>rep f m -&gt;</pre>	
<pre>rep f' m.</pre>	
Proof.	
unfold rep; intuition.	
eapply SDIR.crash_rep; eauto.	
inversion H; intuition subst; cb	n in *.
congruence.	
Qed.	DirCache



	In-Theorem Context
	Additional Context
.V	

### Proofs may be **similar to other proofs**

```
Theorem crash_rep : forall f f' m,
    BFILE.file crash f f' ->
    rep f m ->
    rep f' m.
  Proof.
    unfold rep; intuition.
    eapply SDIR.crash_rep; eauto.
    inversion H; intuition subst; cbn in *.
    congruence.
                                      DirCache.v
  Qed.
```

Theorem crash_rep : forall f f' m,	
BFILE.file_crash f f' ->	
rep f m $->$	
rep f' m.	
Proof.	
unfold rep; intros.	
repeat deex.	
eexists; intuition eauto.	
<pre>eapply DIR.crash_rep; eauto.</pre>	
Qed.	DirName.v

### Proofs may be **similar to other proofs**

```
Theorem crash rep : forall f f' m,
    BFILE.file crash f f' ->
    rep f m ->
    rep f' m.
  Proof.
    unfold rep; intuition.
    eapply SDIR.crash_rep; eauto.
    inversion H; intuition subst; cbn in *.
    congruence.
                                      DirCache.
  Qed.
```



	Theorem crash_rep : forall f f' m,	
	BFILE.file_crash f f' ->	
	rep f m ->	
	rep f' m.	
	Proof.	
	unfold rep; intros.	
$\neg$	repeat deex.	
	eexists; intuition eauto.	
	<pre>eapply DIR.crash_rep; eauto.</pre>	
V	Qed.	DirName.v

## What context should be provided?

Theorem crash_rep : forall f f'	m,
BFILE.file_crash f f' ->	
rep f m $->$	
rep f' m.	
Proof.	<b>†</b> Prompt
	↓To Generate

### Vanilla

Provide **statements** of all accessible definitions, theorems, and lemmas

### Hint

Additionally provide **proofs** of 50% of accessible theorems and lemmas, selected at random

## Evaluation

- 1. How much of FSCQ can be proved?
- 2. What is the effect of proof context?
- 3. Which theorems in FSCQ are easier (or harder) to prove?
- 4. Are models simply memorizing FSCQ's proofs?
- 5. What is the effect of model size?
- 6. When and why do models fail?

## Evaluation

- 1. How much of FSCQ can be proved?
- 2. What is the effect of proof context?
- 3. Which theorems in FSCQ are easier (or harder) to prove?
- 4. Are models simply memorizing FSCQ's proofs?
- 5. What is the effect of model size?
- 6. When and why do models fail?

# Evaluation: How much of FSCQ can be proved?



GPT-40 proves **38%** of the sampled theorems

...and **57%** of the sampled theorems under 64 tokens



# Evaluation: What is the effect of proof context?



Hints improve coverage by between **30%** and **145%** 

Longer context window **does not necessarily** improve coverage





## (With Hints) Expected 32% 42% 47% Based on lengths of human proofs and average proof coverage for the associated range of lengths





	(With Hints)		(Vanilla)	
	Actual Expected		Actual	Expected
File System	20%	32%	16%	24%
CHL	52%	42%	43%	32%
Utilities	58%	47%	40%	36%

### Hints improved performance the most on **Utilities**

## Generated proofs sometimes use lemmas **more efficiently**

```
unfold ndata log, padded log, setlen,
  roundup; intros.
rewrite firstn oob by auto.
repeat rewrite map app.
rewrite repeat map; simpl.
rewrite nonzero_addrs_app.
setoid_rewrite <- app_nil_l at 3.</pre>
rewrite nonzero addrs app zeros; auto
                                  Origina
```

Lemma ndata log padded log : forall a, ndata log (padded log a) = ndata\_log a.

	intros a.	
	unfold ndata_log.	
	rewrite nonzero_addrs_padded_log	•
	reflexivity.	
•		
al		GPT-40

## Generated proofs sometimes use lemmas **more efficiently**

```
unfold ndata log, padded log, setlen,
  roundup; intros.
rewrite firstn oob by auto.
repeat rewrite map app.
rewrite repeat map; simpl.
rewrite nonzero_addrs_app.
setoid_rewrite <- app_nil_l at 3.</pre>
rewrite nonzero_addrs_app_zeros; auto
                                  Origir
```

Lemma ndata log padded log : forall a, ndata log (padded log a) = ndata log a.

	intros a. Additional unfold ndata_log.	Context
	rewrite nonzero_addrs_padded_log	•
	reflexivity.	
•		
nal		GPT-40

### Generated proofs sometimes use lemmas **more efficiently**

```
Lemma tree_name_distinct_head: forall inum name l t,
```

```
intros. destruct t.
constructor. inversion H.
rewrite map_cons in H2.
apply Forall inv in H2.
simpl in H2. inversion H2.
constructor; eauto.
                                  Original
```

- tree names distinct (TreeDir inum ((name, t)::1)) -> tree names distinct t.

```
intros.
inversion H; auto.
inversion H2; subst; auto.
                             Gemini-1.5 Pro
```

# Discussion: Augmenting human effort

## LLMs naturally fit into human-in-the-loop workflows



## They could also complete partial proofs or fill in missing steps

# **Discussion: Reasoning Models**

## What about **reasoning** models?







High resource consumption

# **Discussion:** Proof Decomposition

## Can LLMs generate **intermediate lemmas**?



## **Tactic Generation**

# **Discussion:** Proof Decomposition

## Can LLMs generate intermediate lemmas?

## **Decomposition** Tactic Generation



## Conclusion

LLMs can prove a surprising fraction of FSCQ

Hints significantly improve proof completion rate

Further Discussion / Open Questions

Human Augmentation?

Reasoning Models?

### Can LLMs augment or replace the manual verification of system software?

Proof Decomposition?

**Context Retrieval?** 

**Open Source** (Code and Generated Proofs) 



Gemini 1.5 (128k context, w/ hints)	
Gemini 1.5 (128k context)	
Gemini 1.5 (w/ hints)	
Gemini 1.5	
GPT-4o (w/ hints)	
GPT-40	
0.	.0 0.2 Similar

Generated proofs are **not direct copies** of the originals



