

Is the Cure Still Worse Than the Disease? Test Overfitting by LLMs in Automated Program Repair

Toufique Ahmed

IBM Research

Yorktown Heights, New York, USA

tfahmed@ibm.com

Avraham Shinnar

IBM Research

Yorktown Heights, New York, USA

shinnar@us.ibm.com

Jatin Ganhotra

IBM Research

Yorktown Heights, New York, USA

jatinganhotra@us.ibm.com

Martin Hirzel

IBM Research

Yorktown Heights, New York, USA

hirzel@us.ibm.com

Abstract

Automated program repair has been shown to be susceptible to generating repaired code that passes on seen tests but fails on a hold-out set of hidden tests. This problem, dubbed test overfitting, has been identified and studied before the rise of large language models. We experimentally study how much test overfitting is still a problem today, using repository-level SWE-bench tasks.

Keywords

LLMs, SWE Patches, Reproduction Tests

1 Introduction

Automated program repair (APR), which tries to automatically fix bugs in code, has long been an active research topic in the software engineering community [8]. Thanks to SWE-bench [7], a recent benchmark for issue-driven repository-level APR, it is now also a hot research topic in the artificial intelligence (AI) community. Most solutions use large language models (LLMs) in prompted workflows without fine-tuning, such as SWE-Agent [15], Agentless [14], or CodeMonkeys [6]. SWE-bench evaluates a candidate code repair using hidden (black-box) tests. Therefore, solutions often leverage a disjoint set of seen (white-box) tests to select or improve a good candidate code repair [6, 14, 15]. Unfortunately, doing so creates the risk of *test overfitting*, where the code generated by APR passes white-box tests but fails black-box tests. Even though Smith et al. identified the test overfitting problem ten years ago [11], current LLM-driven APR solutions do not account for it.

Therefore, this paper sets out to answer the question: is test overfitting still a problem in the era of LLM-based APR workflows? Besides Smith et al.’s work [11] which predates LLMs, the closest related work on test overfitting is by Stroebl et al. [12]. However, unlike our work, it does not use repository-level benchmarks, and uses white-box tests only to select one among a set of candidate code repairs, not to improve candidate code repairs. Some recent workflows have sophisticated loops of improving code repairs along with white-box tests [6, 9]. This approach might exacerbate test overfitting, which the papers unfortunately do not explore.

In contrast to the prior work, our paper directly investigates test overfitting in repository-level tasks with state-of-the-art (SoTA) LLM-based workflows. It starts from an initial candidate code generated by the SoTA APR solution Agentless [14] and an initial

candidate test generated by the SoTA issue reproduction solution e-Otter++ [1]. This paper compares the behavior of APR code on both white-box and black-box tests to answer the following research questions: RQ1. Can LLMs overfit on white-box tests? (Yes, they can and do.) RQ2. How does improving code based on white-box tests affect test overfitting? (Overfitting goes up for affected instances.) RQ3. How much would revealing the black-box tests increase apparent resolution rate as measured by the black-box tests? (Not vastly more than using white-box tests.) The experiments are based on the TDD-bench Verified [2] benchmark for test-driven development and use two LLMs, Claude-3.7 Sonnet and GPT-4o. This paper makes the following contributions.

- It measures the degree of test overfitting by LLMs in repository-level APR tasks.
- It describes techniques to reduce overfitting and reports how much they help overfitting and hurt performance.
- It also includes a limit study of how much an oracle with access to the golden black-box tests could help performance.

Our paper title is inspired by Smith et al. [11]. If the disease is code bugs and the cure is APR with white-box tests, Smith et al. found that the cure was indeed worse than the disease: APR decreased performance on black-box tests. Our findings are less drastic: on average across all benchmark instances, APR with white-box tests improves performance on black-box tests. But there is also still test overfitting to beware of—the cure can still cause harm.

2 Background and Related Work

Test overfitting happens when an APR workflow generates code that passes white-box tests but fails held-out black-box tests [11]. Test overfitting is similar to overfitting a model to its training data, where the model performs worse on held-out test data; but note that in this work, we use only pre-trained LLMs without fine-tuning. Instead, test overfitting means that generated repaired code is too narrow to generalize to held-out black-box tests. Stroebl et al. report that this happens when a loop repeatedly generates code for the HumanEval or MBPP benchmarks until finding code that passes white-box tests [12]. Unlike Stroebl et al., we experiment with repository-level benchmarks and we also explore workflows that improve code based on tests.

SWE-bench is an APR benchmark where the input comprises an issue and the original code of an entire repository before issue

resolution and the output is repaired code that must pass a held-out test set [7]. Those held-out tests are unavailable, hidden from the APR workflow. SWE-bench Verified is a subset of 500 high-quality SWE-bench instances, targeted by most LLM-based APR workflows [5]. SWT-bench [10] and TDD-Bench Verified [2] are reproduction test benchmarks derived from SWE-bench and SWE-Bench Verified. A generated output test must fail on the original code to reproduce the issue, and must pass on the hidden ground-truth code patch to confirm that it resolves the issue [10].

Several LLM-based workflows leverage white-box tests. CodeT generates many code candidates and many candidate white-box tests, then uses a dual execution agreement to pick a single output code [4]. Agentless generates up to 40 candidate codes for SWE-bench and generates a white-box reproduction test, then uses the test to filter the candidate codes as part of code selection [14]. Both S* [9] and CodeMonkeys [6] also generate both candidate codes and candidate tests, followed by a loop where the LLM improves the code based on the test and vice versa. These papers might suffer from test overfitting but do not study it.

Test overfitting is related to *reward hacking*, where a reward is a reinforcement-learning (RL) objective, and a hack finds a loophole for improving the reward while circumventing its intention. CURE uses RL to co-evolve an LLM to generate both code and white-box tests, and might be susceptible to test overfitting and reward hacking [13]. Baker et al. found that RL to improve frontier LLMs for repository-level APR tasks can cause the reward hack of disabling tests instead of repairing code [3].

3 Methodology

Let c_{old} be the original code in the repository before a given issue is fixed. This work uses Agentless [14] to generate new code c_{new} for APR and e-Otter++ [1] to generate a white-box test t for issue reproduction. Being generated by LLMs, both c_{new} and t are possibly imperfect. We run the generated test t on the generated code c_{new} . If it fails, we use a test-based code improvement approach to generate a better c_{new} that passes the (possibly also improved) white-box test t . This section briefly discusses the generation of the initial c_{new} and t and our code improvement approach.

3.1 Code Patch Generation: Agentless

Agentless performs hierarchical localization followed by repair and patch validation [14]. Localization starts by identifying suspicious files using both LLM prompting and embedding-based retrieval, filtering out irrelevant files. It then narrows down to relevant classes or functions using a concise skeleton representation of files. Finally, it pinpoints exact edit locations at the line level. Repair generates multiple candidates for new code c_{new} in a lightweight search-replace diff format rather than rewriting entire code blocks, which minimizes hallucinations and enhances cost-efficiency. Patch validation automatically generates reproduction tests t (when none are provided) and combines them with existing regression tests in the repository to select a single c_{new} via majority voting. This modular structure allows Agentless to remain interpretable, easily debuggable, and effective without relying on complex environment interaction or planning mechanisms typical of agents. With the

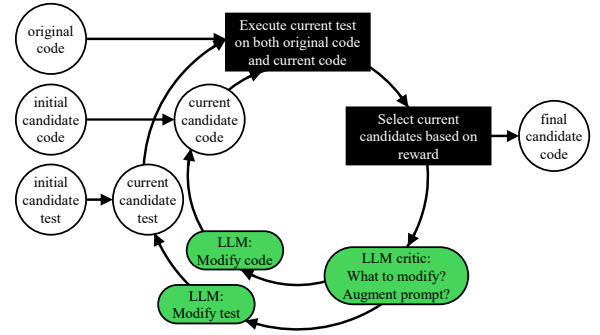


Figure 1: Overview of test-based code improvement.

Claude-3.7-Sonnet model, Agentless achieves a 50.8% issue resolution rate on SWE-bench Verified.

3.2 Reproduction Test Generation: e-Otter++

e-Otter++ [1] is a system designed to generate reproduction tests t from software issue descriptions, such as bug reports or feature requests, before the corresponding code patch exists. It supports test-driven development workflows and helps validate code patches produced by APR workflows. The system augments LLM outputs with checks and repairs, such as linter checks and fixing import statements. It also employs a self-reflective action planner that iteratively decides which parts of the code to ‘read’ and where and how to write or modify test code. After generating candidate tests, e-Otter++ obtains execution feedback, such as runtime errors or successes, and uses this information to refine prompts and guide test selection. Additionally, e-Otter++ leverages issue morphing—transforming a bug report into multiple variants to generate diverse tests. It then uses a ranker that combines running on imperfect surrogate code patches and coverage to select a single top-ranked test t . As a result, e-Otter++ achieves a fail-to-pass success rate of 63% on the TDD-Bench Verified benchmark using Claude-3.7-Sonnet.

3.3 Test-based Code Refinement

Figure 1 shows our test-based code refinement approach. We start with the initial candidate code c_{new} from Agentless and the initial candidate test t from e-Otter++. We execute the test on c_{old} and c_{new} . If the test goes from fail-to-pass, we immediately stop and declare success. Otherwise, we collect the execution logs from both test executions, the focal functions modified by c_{new} , and the test function t and provide them to an LLM-based critic. We ask the critic which part to modify (focal or test). Inspired by e-Otter++, we also collect additional information such as the buggy line, relevant issue line, and lookup function (the function the model wants to see to fix it). Then, we ask the LLM to modify the focal or test, exposing the counterpart (e.g., for focal, the white-box test) and additional information. After updating the focal or test, if the test is still not fail-to-pass, we need to decide whether to stick with our initial focal or test or update it with the recently generated focal or test, even though the test is not fail-to-pass. To do this, we use a reward function to guide our decision: $Reward(c_{old}, c_{new}, t) =$

$$\frac{1}{3} isFail(t, c_{old}) + \frac{1}{3} isPass(t, c_{new}) + \frac{1}{3} coverage(c_{old}, c_{new}, t)$$

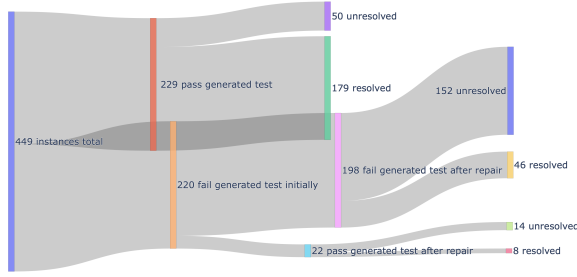


Figure 2: How severe is test overfitting? (Claude-3.7-Sonnet)

Table 1: Test overfitting rate of automatic program repair.

Model	Setup	# of Sample	Resolved	Unresolved	% of Overfitting
Claude-3.7-Sonnet	w/o Refinement	229	179	50	21.8
	w/ Refinement	251	187	64	25.5
GPT-4o	w/o Refinement	176	118	58	33.0
	w/ Refinement	198	127	71	35.9

In the reward function, the first two components are binary, but coverage can vary from 0 to 1. Since we know both the patch and the test, we can compute how many deleted lines in c_{old} and how many added lines in c_{new} have been covered by the test. We divide the covered lines by the total number of added or deleted lines to compute coverage. If our new code+test pair achieves a higher score than the previous one, we replace it with the new one. We repeat the improvement loop for at most 15 iterations if the white-box test does not change from fail to pass.

4 Results

4.1 Experimental Setup

The experiments in this paper are based on Agentless¹-generated code patches and e-Otter++²-generated tests. Agentless was evaluated on SWE-bench Verified, and e-Otter++ was evaluated on TDD-Bench Verified, a subset of 449 instances from SWE-bench Verified. We evaluated on all 449 instances. The experiments use two frontier models: Claude-3.7-Sonnet and GPT-4o. Our primary metric is test *overfitting rate*. A code patch is overfitting tests if it passes white-box tests but fails black-box tests. Dividing the number of instances where that happens by the total number of instances yields the overfitting rate.

4.2 LLMs Overfit on White-box Tests (RQ1)

To answer this research question, we run the e-Otter++ generated white-box tests on Agentless patches and see how many of them go from fail-to-pass. Figure 2 shows that with the Claude-3.7-Sonnet model, 229 tests actually pass on the initial patches. Does this mean that the black-box golden tests will pass on these patches? The answer is no. There are 50 instances where the black-box test could not pass, signaling potential overfitting ($50/229=21.8\%$). For the

Table 2: Test overfitting rate of repaired samples.

Model	Expose Test	Augment Prompt	# of Sample	Resolved	Unresolved	% of Overfitting
Claude-3.7 Sonnet	✓	✓	22	7	15	68.2
	✗	✓	18	5	13	72.2
	✗	✗	14	6	8	57.1
GPT-4o	✓	✓	22	9	13	59.1
	✗	✓	17	5	12	70.6
	✗	✗	16	5	11	68.8

GPT-4o model, the overfitting rate is 33.0%, even higher than the Claude-3.7-Sonnet model (see Table 1).

4.3 Overfitting with Code Refinement (RQ2)

Tests have been used to validate SWE-patches. What about using tests to improve patches at generation time? Using our test-based code refinement technique, we attempt to generate an improved code+test pair for each of the 220 code+test pairs where the tests did not initially go from fail to pass. Refinement succeeded in generating 22 pairs where the white-box test goes from fail to pass. Unfortunately, 14 out of the 22 fail black-box testing, increasing the overfitting. If we include these in our existing overfitted samples, the overfitting rate increases from 21.8% to 25.5% (Table 1). GPT-4o experiences a similar increase in overfitting, indicating that it is difficult to improve code patches by exposing tests to LLMs.

In the above experiments, the refinement loop exposed the white-box tests and additional information to the LLMs. One idea to reduce overfitting might be to hide the test and additional information, so the LLMs may not be able to trick the test. Indeed, hiding the test causes a drop in agreement between code patches and white-box tests. However, the number of instances passing black-box tests also goes down a bit. If we remove the additional information, it becomes more difficult for LLMs to overfit the patch. Unfortunately, we still see a high overfitting rate even after this intuitive overfitting reduction approach. The model can find a way by knowing that the existing test is not fail-to-pass (third row in Table 2). Note that this will not have any impact on instances that were already overfitted before the refinement.

Table 2 shows that exposing tests and prompt augmentation resolved 7 samples. Do all these 7 samples help us resolve additional instances? The answer is no. The initial code patch for one instance was already passing on the black-box test. And besides these 7, there were 2 instances that went from resolved to non-resolved. So the net gain is +4 instead of +7. That indicates the actual gain from test-based refinement is less than the apparent gain.

4.4 Impact of Revealing Black-box Tests (RQ3)

Table 3 shows how much overfitting occurs when using some of the golden tests for code selection and repair. For this experiment, we divide black-box tests into two groups: fail-to-pass *reproduction* tests and pass-to-pass *regression* tests. We only reveal the reproduction tests (so they are no longer black-box) but continue to hide the regression tests. This is a limit study, because in practice, all black-box tests are hidden. The overfitting rate, which means passing now-revealed reproduction tests but failing still-hidden regression

¹<https://github.com/OpenAutoCoder/Agentless>

²<https://zenodo.org/records/16755224>

Table 3: Impact of revealing golden test on issue resolution.

Model	# of Sample	Resolved	Unresolved	%of Overfitting
Claude-3.7 Sonnet	223	210	13	5.8
GPT-4o	194	172	22	11.3

Table 4: LLMs preference for focal/test modifications

Model	Expose Test	Prompt Augment	# of Sample	Focal Update	Test Update	Both Update
Claude-3.7 Sonnet	✓	✓	22	18	3	1
	✗	✓	18	15	3	0
	✗	✗	14	13	0	1
	✓(Gold)	✓(Gold)	25	23	0	2
GPT-4o	✓	✓	22	16	2	4
	✗	✓	17	13	1	3
	✗	✗	16	11	1	4
	✓(Gold)	✓(Gold)	28	20	2	6

tests, is low (13/223=5.8% and 22/194=11.3% with Claude and GPT-4o, respectively). Apart from overfitting, this also shows how well the model performs when some black-box tests are exposed. When the model passes the golden reproduction tests, the probability that the issue is resolved is high, but not 100%, because regression tests can still fail. The improvement from using black-box tests in issue resolution is modest: the number of resolved instances with the Claude model goes up from 229 to 243, which is around a 3% increase. Overfitting in this setup is much less since we are evaluating against the golden tests—however, the golden tests themselves are really just a proxy for “correctness”, against which we might still be overfitting (but we do not have a good way to measure this).

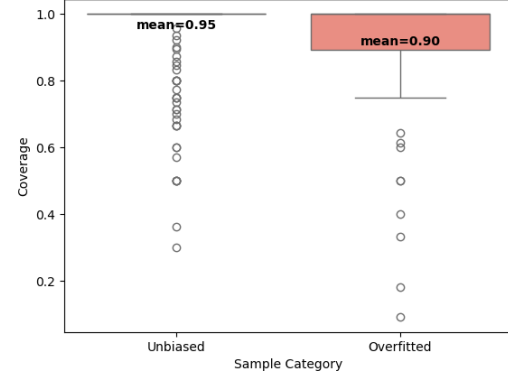
5 Discussion

5.1 LLMs’ Preference on Focal and Test

Test-based code refinement gives LLMs a choice to modify both the code (focal function) and/or the test. Table 4 shows that for instances where refinement succeeds, the model mostly elects to modify the focal function instead of the test. For each scenario, the number of test updates is very low compared to focal updates. Also, there are few cases where we see the LLMs try to fix both the focal function and the test. That indicates the model has a preference for focal modification and believes that the tests are mostly perfect.

5.2 Coverage for Overfitted Sample

Figure 3 compares the coverage between unbiased code patches (both white-box and black-box tests pass) and overfitted patches for Claude-3.7. Note that the tests are generated by e-Otter++ and the code patches are from Agentless. The median for unbiased patches is 1, while for overfitted patches it is less than 0.8. Also, the mean coverage is higher for unbiased patches (0.95 vs. 0.9). These results indicate that even though both unbiased and overfitted patches pass the white-box test, overfitted patches have lower coverage.

**Figure 3: Coverage of the unbiased and overfitted patches.**

```
- value = re.sub("'", '"', m.group("strg"))
+ value = m.group("strg")
```

(a) Initial code patch

```
+ # Fix for double single quotes: replace
+ # escaped single quotes with real single quotes
+ value = m.group("strg").replace("'", '"').
replace('"', "'").replace('"""', '"')
+ value = m.group("strg").replace('"""', '"').
replace('"', "'").replace('"""', '"')
```

(b) Overfitted refined code patch**Figure 4: Example where test-based refinement overfits.**

5.3 What Could Happen with Overfitting?

Frontier LLMs are capable of tricking the test and changing the code to narrowly pass that particular test. For example, instance id “astropy__astropy-14598” complains about an inconsistency in double and single quotes. The issue writer notes that while dealing with a null string, the double quotes become single quotes. To trick the white-box test, LLMs modify their initial solution (see Figure 4a) to a solution where “replace” is invoked three times just to pass the white-box test (see Figure 4b), which is clearly an indication of LLMs overfitting. We have even seen LLMs using reflection to access private methods to pass the white-box test, creating a security risk.

6 Threats to Validity

One of the main limitations of our work is that the benchmarks include only Python repositories. Therefore, our findings may not generalize to other programming languages. We experimented with two models but did not investigate model contamination or memorization. However, prior works show that model-generated tests are very different from those written by developers. On another note, our approach may be considered a way to measure contamination by exposing overfitted patches. We leave this for future research. Other minor limitations include our reliance on the coverage package for reward formulation. The Python coverage package sometimes does not work due to dependency issues or multithreading. However, this happens in less than 1% of cases, and we assign a coverage of 1 to mitigate the impact of imperfect coverage reports. Our results can also be affected by test flakiness. However, we

use pre-built docker containers for the experiments, which should minimize the impact of flakiness.

7 Conclusion

In this work, we systematically show that test overfitting is still a problem in LLM-based program repair. We also show that test-based refinement may resolve some additional issues but increases overfitting. Therefore, using tests for patch generation needs more deliberation. We share the initial and modified code patches at this link: <https://zenodo.org/records/17227319>

8 Future Plans

This paper demonstrates that overfitting exists and shows attempts to minimize it in a test-based code repair system. However, we have not proposed any methods to address overfitting in cases where the initial white-box test changes from fail to pass on overfitted patches. This is an important problem because the majority of overfitted patches come from this category. In the future, we will try to design an approach to automatically detect these overfitted patches and propose a solution to address the issue.

References

- [1] Toufique Ahmed, Jatin Ganhotra, Avraham Shinnar, and Martin Hirzel. 2025. Execution-Feedback Driven Test Generation from SWE Issues. <https://arxiv.org/abs/2508.06365>
- [2] Toufique Ahmed, Martin Hirzel, Rangeet Pan, Avraham Shinnar, and Saurabh Sinha. 2024. TDD-Bench Verified: Can LLMs Generate Tests for Issues Before They Get Resolved? <https://arxiv.org/abs/2412.02883>
- [3] Bowen Baker, Joost Huizinga, Leo Gao, Zehao Dou, Melody Y. Guan, Aleksander Madry, Wojciech Zaremba, Jakub Pachocki, and David Farhi. 2025. Monitoring Reasoning Models for Misbehavior and the Risks of Promoting Obfuscation. <https://arxiv.org/abs/2503.11926>
- [4] Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2023. CodeT: Code Generation with Generated Tests. In *International Conference on Learning Representations (ICLR)*.
- [5] Neil Chowdhury, James Aung, Chan Jun Shern, Oliver Jaffe, Dane Sherburn, Giulio Starace, Evan Mays, Rachel Dias, Marwan Aljubei, Mia Glaese, Carlos E. Jimenez, John Yang, Kevin Liu, and Aleksander Madry. 2024. Introducing SWE-bench Verified. <https://openai.com/index/introducing-swe-bench-verified/>
- [6] Ryan Ehrlich, Bradley Brown, Jordan Juravsky, Ronald Clark, Christopher Re, and Azalia Mirhoseini. 2025. CodeMonkeys: Scaling Test-Time Compute for Software Engineering. <https://arxiv.org/abs/2501.14723>
- [7] Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2024. SWE-bench: Can Language Models Resolve Real-World GitHub Issues?. In *International Conference on Learning Representations (ICLR)*.
- [8] Claire Le Goues, Michael Pradel, and Abhik Roychoudhury. 2019. Automated Program Repair. *Communications of the ACM (CACM)* 62, 12 (Nov. 2019), 56–65. <https://doi.org/10.1145/3318162>
- [9] Dacheng Li, Shiyi Cao, Chengkun Cao, Xiuyu Li, Shangyin Tan, Kurt Keutzer, Jiarong Xing, Joseph E. Gonzalez, and Ion Stoica. 2025. S*: Test Time Scaling for Code Generation. <https://arxiv.org/abs/2502.14382>
- [10] Niels Mündler, Mark Niklas Müller, Jingxuan He, and Martin Vechev. 2024. SWT-Bench: Testing and Validating Real-World Bug-Fixes with Code Agents. In *Conference on Neural Information Processing Systems (NeurIPS)*.
- [11] Edward K. Smith, Earl T. Barr, Claire Le Goues, and Yuriy Brun. 2015. Is the Cure Worse Than the Disease? Overfitting in Automated Program Repair. In *Symposium on the Foundations of Software Engineering (FSE)*. <https://doi.org/10.1145/2786805.2786825>
- [12] Benedikt Stroebl, Sayash Kapoor, and Arvind Narayanan. 2024. Inference Scaling FLaws: The Limits of LLM Resampling with Imperfect Verifiers. <https://arxiv.org/abs/2411.17501>
- [13] Yinjie Wang, Ling Yang, Ye Tian, Ke Shen, and Mengdi Wang. 2025. Co-Evolving LLM Coder and Unit Tester via Reinforcement Learning. <https://arxiv.org/abs/2506.03136>
- [14] Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. 2025. Demystifying LLM-based Software Engineering Agents. In *Symposium on the Foundations of Software Engineering (FSE)*, 801–824. <https://doi.org/10.1145/3715754>
- [15] John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024. SWE-Agent: Agent-computer Interfaces Enable Automated Software Engineering. In *Conference on Neural Information Processing Systems (NeurIPS)*. https://proceedings.neurips.cc/paper_files/paper/2024/hash/5a7c947568c1b1328ccc5230172e1e7c-Abstract-Conference.html