# CS203 Gene's Class - Honeywords Overview

**LEC3-honey-passwords.pdf**
PDF

read this presentation slides about honey password. read and. gather all your previous knowledge

Alright, I've read through your uploaded slides on **honey passwords** ("LEC3-honey-passwords.pdf") and combined it with my prior knowledge.

Here's a **gathered summary** — in a clean, intuitive, complete way, merging the lecture and broader context:

---

## 🧠 What are Honey Passwords (Honeywords)?

- **Honeywords** are **decoy passwords** — fake passwords stored alongside a user's real password.
- If an attacker steals the password file and tries a honeyword instead of the real password, the system can detect the breach.

Honeywords help **make password cracking detectable**, not necessarily preventable.

---

## 🏰 Motivation: Why Honeywords?

- Systems can **hash passwords** with salting and slow-hashing (e.g., bcrypt), but...
  - **Real user passwords** are weak (short, guessable, reused).
  - **Password crackers** model human behavior (using leaked password datasets like RockYou).
  - Even salted and hardened hashes **can still be cracked** over time.

So, we **assume compromise is inevitable** — **how to detect it** becomes the key question.

---

## 🎯 How Honeywords Work: The Core Game

1. **Instead of storing one password per user**, the system stores a **list** of sweetwords (one real password + multiple fake ones).

   Example:

   ```go
   P1, P2, ..., Pi (real password), ..., Pn
   ```

2. During login:

- The password attempt is **matched** against all sweetwords.
- Only the **Honeychecker** (a separate, secure system) knows which index `i` corresponds to the real password.

3. **If a fake password is used**:
   - The Honeychecker raises an **alarm**.
   - Breach detected!

## ➡️ Important Properties:

- **No noticeable change** in user experience.
- **Distributed security**: compromise of either the main system or Honeychecker alone isn't catastrophic.
- **Honeychecker can be minimal or offline**, waking up only to report alarms.

---

# 🛠️ How are Honeywords Generated?

This is **tricky** because honeywords must look real enough to fool an attacker!

**Two main strategies:**

1. **Chaffing with a Password Model:**
   - Use leaked password datasets (like RockYou) to train a model on human password habits.
   - Generate honeywords from the model so they "look real."

2. **Chaffing by Tweaking:**
   - Based on the observation that people modify passwords slightly during resets (e.g., HardPassword1 → HardPassword2).
   - Slight variations of the real password are generated to create believable honeywords.

## ⚠️ Problem:

- Some passwords (e.g., "Blink-182") have **semantic meaning**.
- Simple tweaking would break the logic (e.g., "Blink-183" looks wrong).
- Future challenge: **natural language processing** techniques to handle such passwords intelligently.

---

# 🔥 Advantages of Honeywords:

- **Fast, lightweight** to deploy.
- **Distributed defense** with minimal modification to existing systems.
- **Raises alarms early** before widespread exploitation.
- **No impact** on normal login behavior (unless an attacker is trying passwords).

---

# 🚨 Attack Scenarios and Limitations

1. **Denial-of-Service Risk:**
   - A malicious adversary could intentionally use honeywords for many users, locking everyone out and causing chaos.

2. **Under Duress Authentication:**
   - Could extend the system where entering a particular honeyword signals a distress condition (e.g., user under threat).

# ✨ Overall Takeaways

- Honeywords don't *prevent* password theft — they **detect** it early.

- They are a practical application of **deception in security**.

- The larger research goal is to **turn deception from art into science**.

Would you also like me to draw a **diagram** of the Honeywords system (user → main computer system → honeychecker flow)?
It might make it even easier to visualize! 📈✨

Also, I can create a **real-world example scenario** if you want to see how it would actually play out. Want that?

> explain it
> Given ideal honeywords, the
> attacker will guess correctly, j = i,
> with probability about 1/n!

Good question — and it's subtle. Let's **break it down intuitively**:

## The situation:

- Suppose you store **n sweetwords** for Alice:

  ```arduino
  [P1, P2, ..., Pi (true password), ..., Pn]
  ```

- An attacker **steals the whole password database**.

- The attacker **doesn't know** which one is the real password — Pi — and which ones are honeywords (fake).

## What happens when they try to impersonate Alice?

They **have to guess**:
"Which one is the real password among the n sweetwords?"

## Now, the key line you asked about:

"Given ideal honeywords, the attacker will guess correctly, j = i, with probability about **1/n**."

- "**j = i**" means: attacker guesses password number `j`, and it's exactly the correct one `i` (true password).

- **1/n probability**: If the attacker has no better information (honeywords are **indistinguishable from real ones** — *ideal honeywords*), then **random guessing** is their only option.

- So **probability of success** =

$$\frac{1}{n}$$

where `n` = total number of sweetwords.

## ⚡ Example:

Suppose there are 10 sweetwords:

```csharp
[P1, P2, P3, ..., P10]
```

- The real password is P4.
- Attacker steals the list.
- No way to tell which is real because they all look legit.
- So attacker picks randomly: **10 options → 1/10 = 10% chance** of getting it right.

## 📌 Quick Note:

- Your slides wrote "**1/n!**" ("n factorial") **by mistake** there —
  it **should be "1/n"**, not "1/n!".

- **n! (n factorial)** is for permutations (like ordering n items), but here it's just picking one correct item out of n → 1/n.

So:

| Slides | Correction |
|--------|------------|
| 1/n! | Should be 1/n |

✅ **Corrected statement**:
*"Given ideal honeywords, the attacker will guess correctly, j = i, with probability about 1/n."*

## Final Summary (One-liner):

👉 If honeywords are well-designed and look real, an attacker who steals the password database has only a **1/n chance** of guessing the correct password for a user.

Would you also want a visual/flowchart showing the attacker's decision process? (Might make it even easier to picture!) 🎯

this is the transcribed class recording on this topic, read it very carefully, understand and keep it in your memory

Alice's password black, it's determined. Well, that means the adversary commits to a certain J, right? Where J is a P, J, right? Obviously, the adversary knows all the passwords already, right? All N of them, he already knows in the clear.

But which one is really Alice's password? The probability of guessing is 1 over N, correct, guessing correctly.

So, there she goes, there he goes, he says, I'm Alice, here's a password P, J, indeed P, J exists, right? And P, J will be authenticated correctly by the server, right? So, these are all passwords.

One of them is Alice's, and one of them the adversary will enter. Maybe the same as Alice's, but yeah, with probability 1 over N.

Which one is Alice's password? Can you tell me by looking at this, what could be her password? This is not a trick question. There's no, there's nothing hiding in there. It's just a set of random cram.

So, but we're not done, right? So, the adversary will likely make a mistake, right? By guessing an incorrect password.

But what now? So, in other words, how do you verify that it's real Alice or an intruder, an attacker? Well, this is where the cute little thing comes in, and I'll show you in a minute. That's one question, right?

So, how does the system, right? The system verify that this J is from the adversary, this password J is from the adversary, but password I is from Alice?

And the second, how do we generate these hiding words? So, they all look like what you saw earlier, kind of like indistinguishable or equally likely to be a password.

And you can see that if you generate random passwords that sound like, that look like no human could remember them, and one of the passwords for Alice is something predictable like my friend Bob, then clearly that is the password.

So, that's a tricky question, right? How do we make the bogus passwords, the decoys look real? And then there's, okay, how do you respond?

So, the idea is very simple. You have the system that has, I mean, I'm just showing this as a, the system that has the password, the fc password file.

And it has basically all this, it has means of verifying all these passwords. But in addition to the system that exists today, okay, this is what we already have today.

But instead, I mean, in addition to that, we also introduced a new component, physically separate component, called the honey checker.

And a honey checker's only task is to remember, for every user, the index. The honey checker knows no passwords.

So, think of it as like a really, like, brain damage component, like really stupid component. It just does one thing.

It gets a user name and an index. And it says, yes or no. Yeah? That's it. That's pretty much 99% of it.

So, when Alice logs in, right, and supplies password i, the system will check the password i as it

already does today.

Nothing changes so far. But in addition, once they check fast, right, once the password i has been authenticated as Alice's password,

the index i is sent over to the honey checker. And the honey checker does a quick lookup in its little database and says,

Alice has index 5. Is i 5? Yes. All good. Nothing happens. But if i is anything other than 5, well.

So, Alice applies the password. The password is matched to pi, right?

You're so far with me, right? It's not exactly that easy because actually the honey checker needs to compute,

remember, this assaulted hash? But he doesn't know what the i is, right?

So, remember, over here, over here, the password, I show them as being here, but they're not stored like this, right?

Do you remember how they are stored? They are stored in the hash, salted hash form.

So, when Alice supplies the p, the system needs to look up all the Alice's, all the m-t-gerts-part into Alice's,

on each one, compute the salted hash, right? Because the salts are different. So, she computes the salted hash,

then she says, oh, there is a match to pi. And then it sends i over to the honey checker.

And the honey checker just says, yes, does nothing. Or, if the adversary or somebody supplies,

well, it has to be the adversary, actually, if they combine it, be j. And the j, and at this point, the system does not know.

Notice, this big box here, this system has no idea what the index for Alice is, what the index for Bob is.

It does not know any indices. It just knows it has m entries for every user.

So, all it does, it authenticates the user using the supplied password. If authentication succeeds, it says, oh, this corresponds to j,

j, and it sends j to the honey checker. Bam. Wrong. Aligned. So, now you see the whole thing.

Most of the time, the honey checker does absolutely nothing other than just do a lookup. If a lookup succeeds, all good.

That means that j matches i. But if a j is not i, alarm. What does it actually mean?

Notice, if somebody is trying to guess Alice's password, if the adversary does not have the password file,

but it's trying to guess Alice's password. With high probability, with all-knowing probability,

the password that the adversary guesses will not hash into any of Alice's honey passwords, right?

There's only an N of them. So, that's life as usual. The system will say, wrong password, as it does today, right?

And you'll say, just wrong password. Invalid, invalid username, password combination. Yeah? Okay.

So, that experience doesn't change. Alice mistypes her own password.

Forget to shift, cap, whatever. No problem. Wrong username, password combination. Try again.

Same as you do today. But the only way, the only way, again, with high probability, the only way that somebody will guess, or have in their possession,

Alice's password, which is not hers, right? Which is one of the, not i's, right? Anything but i, is going to be somebody who cracked the password file.

Does that make sense? That's the cause for alarm.

So, you wouldn't want any of these alternate passwords to be a very commonly known password?

Right. Right. You don't want them because then you'll be confused, right? Is it, because it has, it cannot be something that somebody would guess about Alice, like trivially, right?

So, that's the, that goes back to how do you generate these details, right? The other question. How do you generate? We are not there yet.

So far, so good? So, yeah, this is already what I mentioned. If the true password is submitted, the user is authenticated.

If a password that is not in the set of Alice's password is submitted, that's a normal authentication error. Okay?

If, if a non Alice's, if a password is submitted, that is one of Alice's, but is indexed, not i, then an alarm is raised.

Only a breach, right? Can cause this. Now you should be convinced. And as you said, they have to be non-treated.

So, if you choose the honey words, judiciously, they will not be like a, a coincidence, right? That you, that somebody, without breaking the password file,

typed in one of the honey passwords. But interestingly enough, as far as Alice, Bob, Charlie, any user is concerned,

there is no change in their experience, right? They don't need to remember two passwords. They still remember one.

They still enter user name, password, hit return. Nothing changes. That's another nice thing.

You don't want to change the user experience, right? Because user experience requires change, requires training. Adaption. Right?

Different users adapt differently. But if you tell them, oh, it's the same as before, just there is something in the back.

Maybe you don't even tell them. They don't need to know.

So, the nice feature here is that the system that today hosts the password file that authenticates you, just needs to transmit the index.

Like minimal information, right? Just name of the person, of the user, and the index.

And so, very little modification, right? You just introduce kind of a call, right? An additional call after you have authenticated the password, right?

Here. After the password has authenticated, you take the index, and the user name, and send it to the honey checker.

And you wait for reply, right? And if the honey checker says yes, cool, honey checker says no, alarm, lock up everything, yeah.

And this is kind of a very trivial way of implementing distributed security, right?

So, the honey checker, of course, is not a full-blown computer system with accounts, right?

You want the honey checker to be a device or a computer that is at most as one account, I mean, basically, like, administering,

user, right? Or a root account that is touched and logged in almost never, right?

Only when the password file and indices are updated.

In fact, you can keep indices the same even if you change the passwords, because they are only known to the honey checker.

No single point of compromise. So, for example, if the computer system is compromised, well, we said, right,

the password file is a hack, you learn about it, right? Because the first user, the first time the adversary

threshold again for any user, for any Alice, Bob, and Charlie, alarm will be erased, and that's it.

The adversary gets one chance. And so, whatever user he picks, Bob, Charlie, Alice, assuming, let's say,

n is some reasonable number, like 20, his chances of getting in a one out of 20.

He might get like, but he's got only one chance, and it's one out of 20.

How long would you update the password without knowing you being fixed?

Well, that has to be done offline. That's not the pretty part. Because if you want to update, if you, well, let's think about it.

You can do it if they gave you their old password.

Yeah, yeah, yeah. I mean, password change, password change programs are separate. Actually, that's a whole, like, separate lecture.

Doing password changes is an underappreciated problem, and a difficult one. Especially considering

what happens if you have disconnects in the process.

Meaning that, at some point, the protocol for changing password breaks, like because of the system malfunction or network disconnect,

and you have essentially a user who thinks the password changed, the system does not, or the other way around.

So, that's a careful thing. That's a problem with all password changes, right?

So, the only way to assure your password changes is to actually manually log in into the system and password change,

instead of remotely, physically close. All right? But, that's not scalable, right? We don't want to do that.

So, that's a whole different headache that is not really different, so much different from the, the case when you use Honey,

Honey passwords. But, let's see. If the Honey checker fails, right? It just, the system still will function, right?

Like, if Honey checker goes down, it uses a hardware failure, or some sort, or it gets disconnected, the system will function. It will just default to what we have today. It will get no worse than what we have today.

So, if the computer system fails, well, you can't log in, but that's the same experience we have today, right?

So, nothing changes so much, if either component, right, is compromised, you don't, it's not failed.

If you, well, sorry, if it fails, if it fails, it's not failed. Now, if either is compromised, like I said, if the computer system is compromised,

well, that's the case we already assumed, right? That the adversary compromised it, and learned the password.

Now, if somebody compromised the Honey checker, without compromising the system, right? So, okay, this is secure,

the adversary does not know the et cetera password file. It just breaks in here, oh, the adversary learns usernames and indices,

but that's it. That doesn't help the adversary to log in. Now, what happens if the adversary compromises both, which is, of course, should be unlikely.

If the adversary compromises both, it gets no worse than we already have.

If you think it through, the situation is strictly no worse than we have today.

Because when the adversary learned the indices, and the adversary learned the password file,

well, he just still has to brute force all the passwords, right, as he would today, so nothing changes.

And Honey checker can be a very, very, like a simple operation.

This may not be, so, if you have, for example, a wire between them, right, like an ethernet wire or some other fiber optic wire

between the computer system and the Honey checker, basically you don't even need any output from the Honey checker.

And the reason is that you know that it received, right? When you send something over a wire, you know it will be received.

Then you have a dedicated wire between the two devices.

So, what if nothing, what if, what if this is a legitimate login?

Well, silence, everything okay?

What if it's illegitimate, meaning that this is an adversary supplying some BJ?

Well, in that case, a physical alarm will go on, right?

The doors will lock, the administrator will be alerted, right?

So, some physical consequences will happen, right?

So, what this is, essentially, it's an input-only system, that you just give it name, index, name, index.

And in case of a problem, it will raise a real alarm.

In practice, you may still want to know as soon as possible, but still.

Make sense?

So, it can be, like, downstream somewhere in an operations center.

So, clearly, you want the Honey checker to be more secure, because this is obviously accessed by, this system will have multiple user accounts, right?

Because they are all logging into it.

They will have the system administrator.

The Honey checker can be, basically, not untouched, for the most part.

So, sitting there, in a box somewhere, does it need a screen, does it need anything, just sit in a box, in a closet, in a, you know, behind an armored door or something, and that's it.

Right?

It just gives this sort of rapid alert.

Right?

So, I reset that.

Okay.

Now, we come back to this issue of Honey word generation.

How do you do it?

Now, look at this.

Which one is the password?

Obviously, the next to last, because it's the only one that looks like human language.

The rest look pretty random.

So, it's very likely, I'm not guaranteed, maybe Alice is a genius, can remember, like one of those savants, can remember a string of two characters, but don't like it.

Right?

So, boom.

Done.

This one.

So, well, what could we do?

Remember, we have Password Crackers, the software, right, that, that, you know, breaks passwords, and it, you know, generates human readable passwords.

So, we can kind of take passwords, right, potentially decoy passwords, from prior breaches of, like, of password databases, and repurpose them.

By adding a letter, a number here, an exclamation, or a question mark, or a comma, or a period, right?

Because they are already human passwords, right?

The ones from, like, a raw human database.

Just grab some, not too popular, but, like, grab some mildly popular password, and tweak it a little bit, and make it a decoy.

Right?

So, let's say it's something like this.

Well, now, now it's difficult, right?

If you're that attacker, you're, you're kind of out of luck, looking at these, these passwords.

They could all be Alice's password.

Happens to be that, but could have been any other?

Okay, there are problem cases.

Now, if you look at this, which would be more likely to be the password here?

My guess is, not wild guess, right?

It's probably the second, right?

It just kind of sticks out at you.

Right?

It's mostly more timely, also.

All right.

So, not good, right?

This is not a good selection.

Now, if every one of them said, I don't know, down with Joe Biden, down with Donald Trump, down with Gavin Newsom, down with, I don't know, Maduro, et cetera.

Yes.

Okay.

That would be hot.

Because they all look very similar.

But, but this, hmm.

So, one way to generate similar believable passwords is by what's called tweaking.

Oh, yeah.

So, essentially, you know, I'll put you to the paper, but of course nobody's going to look at it, but the idea is very similar.

Tweak the actual password.

Take, take the user's actual password that he or she picked and tweak them a little bit to generate, like, minor variations.

That has problems, but, something like this.

So, both lines are passwords, like gamma half, pacificer six, and this is contamination, right?

So, all four of these are passwords.

I mean, nobody, I'm saying, well, he thinks that your password is very long, but, but you get the idea, right?

They all have similar structure.

They all start with, like, the same sort of human readable part, and then there are these numbers, right?

Three, two, one, four, five, six, seven, and then there are the sequential numbers, but they all look very similar.

Now, looking at this, it's very hard to say, right?

So, one of them is, is obviously the password, and the rest are tweaked, but which are tweaked? Why?

Because all the other ones, it has the one that has the most common numbers and letters.

Maybe, maybe, maybe, but it happens to be this one.

No, I mean that, of course, it's a synthetic example, but it happens to be this one.

But you get the idea, chaffing is a, does everybody know what the word chaff means?

Chaffing, like, when you have weak chaff, you know, stuff that is like fluff, noise.

So, the rest of them are noise, right?

So, here's another problem.

Here you have a bunch of very similar passwords, right?

Which one is the real password?

Any guesses?

I know you, you have a guess, but does anybody want to guess?

Okay.

It's the fourth one.

It's the fourth one.

Why?

The band.

The band.

Now, of course, now you know this, except for him, I don't know why he knows, but it's a band from, you know, some, a couple of generations before you.

It's a rock band, right?

So, it has meaning.

The rest of them are kind of random, right?

But this one has a real meaning.

Looks random, but it is not.

So, the semantics of these particular passwords are significant, whereas the rest of them don't have any real semantics.

So, if you just tweak by generation, like if the real password is Blink 182, and you start generating Blink 123, Blink 183, well, none of them are rock band.

So, somebody who knows this will look at it and say, uh-huh, I know exactly what the password is.

That's difficult.

I think today you can get GBT to generate some meaningful passwords.

Maybe, now this stuff wasn't proposed in an LLM era, right?

It was before.

Yeah.

You probably can.

But there's a better idea.

You can just repurpose other people's passwords.

You know, even conceivably on the same system takes Bob's and Charlie's passwords and throw them in his chaff and Alice's mix.

They're as believable as Alice's passwords, you see?

You can do that.

Well, basically, what is this?

What is the idea?

What is the honey password scheme?

It's an example of distributed security, kind of inexpensive distributed security that strictly strengthens the resilience of the system.

Now, it's not going to prevent password breaches.

It has absolutely nothing to do with prevention.

It has everything to do with timely detection.

Okay?

So, it will, it's purpose only one, to detect ASAP when something happens, like a password file is compromised and the adversary tries to gain access.

Yeah?

Can the adversary use these correct passwords and maybe attempt them on a different website?

Of course!

And they're forgetting that they're using the password.

Of course.

Of course.

This is not a panacea.

This is not magic, right?

The adversary stole the password file.

And Greg told his password.

He can try them on a different system because many of us, we use passwords.

We do, unfortunately.

Now, some, some password guideline checkers, you know how when you select your password, you are asked to enter, I mean, must have at least one lowercase, one uppercase, one special character,

and one number, right?

I think usually that's what you see I say.

And the length, right?

Eight characters at least, I think, for us.

What was the last time anybody changed the password?

I think I have.
Okay.
You haven't changed the password either.
I have.
I have.
But I just don't remember the number.
There are these other rules, but I remember, I think it, I think it's eight.
Anyway, so some of the systems will have an additional guideline where you must insert the name of the system you are logging in into your password.
Meaning that this is like uci.edu.
So your password should have uci.edu in it.
That is a pain.
Not because if you have uci.edu, that part isn't useful.
So you have to essentially add this to your password at the end or in the beginning.
You're not going to sprinkle it throughout.
And if the limit is eight, well you should not count these seven characters as part of the password.
And for strength purposes, right?
So essentially, by mandating that you all use uci.edu in your password, I'm saying that your password has to be at least 15 characters long, right?
You see what I'm saying?
Now, but that would solve a problem, right?
In that you could not then reuse, well, it would be awkward for you to reuse that same password because it contains uci.edu on, I don't know, wellsfarga.com.
But that's a balance between password usability and password strength.
Now, back to your question.
There's one other problem.
Yes, the adversary could try.
But remember, the adversary still doesn't know the right index.
So it's very likely it's going to be frustrating for the adversary if the other system also implements a honey checker.
Well, they would have different honey.
That's right.
But then they usually have different passwords.
Right, but they just have to try like 20 times and then one of them works.
No, no, no.
Yes, if the other system, so let's say you broke into, you got uci.edu password for you.
And then you went to Wells Fargo and they tried to log in as u using those passwords.
But if Wells Fargo also implements a honey checker, that ain't going to work, right?
I mean, it ain't going to work.
Actually, no, let me take it back.
It will work because one of them is probably, yeah, not probably, one of them is going to overlap.
So actually, yeah, it might actually work.
You need to compare it.
That's right.
Well, no, no, you don't have to break in to compare.
If you break into Wells Fargo, then you can compare.
But that's a high bar.
I think actually you're right.
You could take them to another website and try 20 times.
Now, after three or four times, you'll get locked out.

So you would have to do this under radar.

Remember I mentioned that a good hacker or adversary does not rush to victory.

If I know 10 passwords for you and only one is real, I'm not going to try one after another.

I'm going to try one and cross it off by myself.

And I'm going to get an education payment.

Then knowing that you log into Wells Fargo once a day, suppose, I will wait a day.

And during that time, you would have logged in successfully.

Okay?

Reset the counter.

I will try the second answer.

Doesn't work?

Cross that off.

Wait another day.

You get what I'm going?

Smart adversary will do that.

Notice that most systems do not warn you when you log in successfully whether they have been in successful attempts before.

Whether they should or not is an interesting question in and of itself.

Because if they start telling you it will be meaningful only if you have good memory, if you log in regularly, you know what I mean?

Like if you are one of those ill retentive people that logs it every day at 8 a.m. to your Wells Fargo account, yes, you will know.

You say, what the hell?

At 4 p.m. last night somebody entered the wrong password.

Ah!

But if you are like most people, you don't remember, you don't have that capacity.

Most of us don't do this already.

Everybody following this conversation?

Yeah?

Okay.

So, anyway, back to this world.

This is a nice balance between deployment and security, right?

Because it's relatively easy to deploy.

Minimally simple, blah, blah, blah.

Now, two follow-up questions and we are done.

That was two interesting things.

Actually, what you raised is also interesting.

I hadn't thought about that.

Can we repurpose this theme to add a separate password?

The Alice would have to remember two, two passwords.

Well, one of them is real password and the other one is I'm being threatened, I'm under duress.

Notice that banks use this system, you know, for PIN codes and combinations for like saves and stuff.

There's two combinations.

One is emergency, one is real.

When you enter an emergency combination, the safe will open.

But they will immediately notify the police or security.

So the idea is the same here.

If you have a list of, if you are Alice and you have a list of N passwords,

it would be nice if you remembered two and only entered that other second one in case of emergency.

Somebody puts a gun to your head and says, you know, log into your bank account and transfer money.
Well, you enter, but the bank now knows that somebody is threatening you.
Yeah, it'd be cool.
And it can be done with the honey checker, right?
The honey checker would raise an alarm when you enter that password that is signified in the emergency.
And now, finally, there is a problem with the honey checker.
Nothing is just beneficial.
One of the things that the industry may aim to just recap it.
Maybe the industry doesn't care to get into the computer system.
He wants to inconvenience you as much as possible.
So you imagine you are a large-ish organization where the system has numerous accounts, maybe thousands, right?
This et cetera password file has like thousands or hundreds of accounts, employees, right?
The adversary breaks one password.
Just one.
Okay?
He wants to take his chances.
He doesn't have to crack all of them, right?
So the answer is easy.
And enters that.
Says, Alice, here's password.
User name Alice, here's password.
Immediately the alarm goes on.
What happens at that point?
At that point, everyone is locked out and everybody must change their password.
That's a giant nightmare.
So the adversary could be very happy just achieving that.
If the adversary's goal is denial of service.
Questions?
No questions.
All right.
All right.
Now it's time to switch to something different.
Let me share the screen.
All right.

Got it — I've **carefully read**, **understood**, and now **integrated** this entire detailed transcribed lecture into memory along with the slides and previous knowledge.

Here's what I **clearly understood and now carry going forward**:

---

## 🧠 Core Concepts Reaffirmed from Lecture:

- **Honeywords** are fake passwords added alongside the real password to confuse an attacker.
- The attacker sees **n passwords** (sweetwords) — one real, (n-1) fake — and **cannot tell** which is real.

- If an attacker guesses:
  - **Chance of correctly guessing = 1/n** (not 1/n! factorial — confirmed).
  - **Probability of mistake = (n-1)/n**, very high when n is reasonable (e.g., 20).
- If the attacker submits a wrong sweetword, the **honeychecker triggers an alarm**.

# 🔒 The Honeychecker: How It Fits

- It's a **separate, physically independent** small device or service.
- It **only stores username and index** (which sweetword is real), **not the passwords** themselves.
- It does **input-only processing**: receives (username, sweetword index) and checks if they match.
- If **yes** → silent, continue normal login.
  If **no** → raise an **alarm** (physical alarms possible like doors locking, alerts firing).

## ✨ Key Points:

- Minimalistic, input-only.
- Very little code or surface area = easy to secure.
- If **honeychecker fails**, system **falls back to standard authentication** — **no worse than today**.

# 🔥 Adversarial Scenarios

- If an attacker **steals password database only**:
  - They still have to guess among n sweetwords.
  - **Wrong guess = alarm**, early breach detection.
- If attacker **steals honeychecker only**:
  - They learn usernames and sweetword indices.
  - **But without passwords**, this is **useless**.
- **Only if attacker compromises both** does security fall back to the old, normal password compromise situation (no worse).

# ⚙️ Honeyword Generation Techniques

1. **Model-based Generation**:
   - Use real password leak datasets (e.g., RockYou).
   - Generate fake passwords that look like believable human-created passwords.
2. **Tweaking**:
   - Slightly modify the real password (increment numbers, slight edits).
   - Mimics how users naturally create password variants.

3. **Semantic Problem:**

- For meaningful passwords (e.g., "Blink-182"), tweaking doesn't work — because semantics give away the real one.

- Solution: Repurpose real passwords from others.

# 🚨 Attack and Limitation

- **Denial-of-Service Attack**:
    - An attacker could intentionally trigger the honeychecker by submitting a wrong sweetword.
    - It would cause **organization-wide password resets**, leading to **chaos**.
    - So even without wanting to log in, an attacker can still do damage.

- **Cross-site Risk**:
    - If password is reused on other sites (and honeywords are not used there), stolen password guesses can succeed elsewhere.
    - **Honeywords don't prevent credential reuse attacks**.

# 🛡️ Key Security Philosophies from the Class

- **Honeywords ≠ prevention**: They **detect** breaches quickly.
- **Distributed security**: No single point of catastrophic failure.
- **Minimal user impact**: No user behavior change.
- **Timely breach detection** is the primary objective.

# 🎯 Final mental picture:

| Component | Role | Danger if compromised |
|---|---|---|
| Main system | Authenticates passwords, stores sweetwords (hashes) | Attackers can trigger honeychecker |
| Honeychecker | Stores user→correct index mapping | Attackers only know mappings, no passwords |

**Attack surface is split — both need to be compromised to defeat the system.**

✅ **All this is now fully stored in my memory** for future discussions.
Would you also like me to now draw an overall **"visual system flow"** for you (how login, honeychecker interaction, and attack detection works)?
It would be super clean to have one! 📈

> did the prof discuss this followup question from slide?
>
> Can honey password scheme be extended to raise an alarm when Alice is
> under duress (e.g., gun to the head)
> – Supplies Pj where j tells the honey-checker that Alice is in trouble

**Yes**, the professor **did discuss** that exact follow-up question from the slide in your transcribed recording — and he explained it quite nicely.
Here's exactly what was discussed:

## 💬 What the Professor Said About Honeywords Under Duress:

- The professor mentioned that **yes**, the honey password scheme **can** be extended to **raise an alarm if Alice is under duress** (e.g., someone forcing her to log in at gunpoint).
- The idea would be:
  - Alice would **memorize two passwords**:
    - One **normal password** (index i).
    - One **duress password** (special index j ≠ i).
- When under duress:
  - Alice would **enter the special password Pj**.
  - The **honeychecker** would recognize that j is the **designated duress index** (different from normal index i).
  - The honeychecker would then **raise an alert** (e.g., notify police/security) **without tipping off the attacker** (because login might still succeed to avoid suspicion).

## 🧠 Key Points the Professor Made:

- This is similar to how **banks and safes** sometimes have **duress PINs**:
  - Entering the duress PIN **still opens** the vault/safe (so attacker thinks all is fine).
  - **Silently alerts** security forces.
- **Technical feasibility**:
  - The honeychecker would need a **mapping of both normal index (i) and duress index (j)** per user.
  - When j is submitted, **special silent alarm** would be raised instead of immediate rejection.
- **User challenge**:
  - Alice now has to **remember two passwords**, not one — a slight increase in user burden.
  - But it could be life-saving in hostage or forced-login situations.
- **Deployment note**:

- This would need a **small modification** to the honeychecker — not just (username, index), but now (username, normal_index, duress_index).

---

# ✅ Final Answer:

- **Yes**, professor discussed it.

- **Yes**, it's feasible.

- **Yes**, it's similar to duress PINs used in banks.

- **Challenge**: Slight increase in user memory burden (need to remember two passwords).

---

Would you like me to also sketch how this **duress honeyword flow** would work in a small simple diagram?
It could help solidify the idea! 🎯 🚨

👤

what did the prof discuss about it?

What if the Adversary's goal is to cause havoc?
– Learns password file, breaks (brute-forces) all users' passwords.
– For each user X, attempts to log in using Pxk – one of X's honey passwords
– Everyone is locked out, must change all passwords ASAP à chaos
– Excellent Denial-of-Service attack!

**Yes**, the professor also **discussed this follow-up question** very thoroughly during the class — and it was quite an important point.

Here's exactly what he said and meant:

---

# 💥 What the Professor Explained About Adversary Causing Havoc (Denial-of-Service Attack):

- The professor said **the adversary might not even want to break into accounts**.

- Instead, **the adversary's goal could be purely to cause chaos** — **a denial-of-service (DoS) attack** using the honeyword system itself.

---

## 🚨 Attack Scenario:

1. **Adversary breaks into the system** and **steals the password file** (the password hashes and sweetwords).

2. **Brute-forces** the file to **recover passwords** for many users (could recover real passwords or the honeywords).

3. Then, **for each user X**, the adversary **deliberately logs in** using a **honeyword** (i.e., wrong sweetword, not the true password).

4. As soon as a **wrong sweetword is submitted**:

   - **Honeychecker immediately raises an alarm**.
   - System **locks the user account**.
   - Organization is **forced to reset passwords for everyone** to prevent further damage.

5. If the adversary does this **for multiple users**, the **entire organization is paralyzed**.

   - **Hundreds or thousands of users locked out**.
   - Huge disruption: **everyone must urgently reset their passwords**.
   - Big **operational nightmare** for IT/security teams.

---

## 📣 Professor's Important Comments:

- This **is a very real, very dangerous attack**.
- The adversary could **achieve denial-of-service** without ever successfully logging in or stealing data.
- **One wrong password submission = massive reaction** (alarm, lockout, forced reset).
- It is a **trade-off** between:
  - **Detecting real breaches early** (which is the purpose of honeywords),
  - and **being vulnerable to chaos if an attacker abuses the alarm system**.
- This attack is **especially dangerous in large organizations** (with thousands of user accounts).
- **He emphasized** that while honeywords improve breach detection, they **introduce new attack surfaces** like this **DoS vector**.

---

## ✅ Final Takeaways:

| Point | Summary |
|---|---|
| Attack type | Adversary deliberately uses wrong honeyword to trigger alarms |
| Goal | Cause massive lockouts, chaos, forced password resets |
| Effect | Denial of Service across organization |
| Trade-off | Detection of breaches vs. vulnerability to abuse |

The professor **clearly warned** about this attack path and highlighted it as **an important limitation** of the honeyword scheme.

Would you also like me to write a **small clean bullet-point list** summarizing both **advantages** and **limitations** of honeywords, as if preparing notes for revision or an exam? 📚✨ It could help you a lot!