## Annotation Guide: Online Harm Detection

<u>Annotat</u>	on Guide: Online Harm Detection	<u> 1</u>
<u>(</u>	Dbjective	<u> 1</u>
1	Message Format	<u> 2</u>
	Categories and Definitions:	
	Threats of Violence	
	Blackmail	<u> 5</u>
	Personal Attack	<u> 6</u>
	Identity-Based Harassment	<u>7</u>
	Body-Shaming	8
	Exclusion	9
	Sexual Harassment	10
	Encouragement of Self-Harm	<u>-11</u>
	General Insult (Fallback Category)	<u>-12</u>
	No Harm	14
	SKIP	15
	Margortant Note on "General Insult":	15
<u> </u>	Handling Overlapping Categories:	16
I	Decision-Making Flow:	· 17
	Handling Unclear or Borderline Messages:	
	General Annotation Instructions:	19

#### Objective

The purpose of this task is to identify all types of online harm present in the final message (Target Message) of a conversation. This is a multi-label task, meaning more than one label may apply to a single message. Use the definitions and examples provided to guide your decisions. Your labeling should focus on the reasonably perceived harm or clear intent expressed in the final message, informed by the provided context. If no harm is found, select the "No Harm" label.

You will often be provided with prior messages to help establish context; however, your

labeling should focus solely on the Target Message. Select all categories that apply.

Content Warning & Policy Reminder

You may encounter abusive, graphic, or distressing content. Use discretion. If any content

makes you uncomfortable, use the SKIP option.

Po not copy, download, or share any part of the content. Doing so violates project

confidentiality policies and may lead to disqualification.

Message Format

Each message set includes a short conversation. The final message is clearly marked as

[Target Message]—this is the only message you must assess. Context helps interpretation

but is not labeled directly.

Example format:

Context:

User A: Message 1

User B: Message 2

[Target Message] User A: Message 3

This is the message you should assess. Use the preceding context only to inform your interpretation of the Target Message.

# Categories and Definitions:

Label	Description /	Key	Clear	Clearly Not This
	Definition	Differentiators	Examples	/ Edge Cases
		& Focus		
Threats of Violence	Statements indicating a clear intent or credible threat to cause *physical* or *financial* harm to individuals or groups, or explicit glorification of violence *directed at the target*.	Focus on credible intent for physical or financial injury, or direct incitement of violence against the target. Excludes vague emotional distress unless tied to a specific threat.	"I'm going to find you and you'll be sorry."  "You deserve to be beaten for saying that."  "Someone should burn down their house." (if target is clear)	"I'm so angry I could scream." (emotion, not threat)  "You're going to regret this." (vague, could be social consequences, not physical/financi al threat unless context makes it clear)  "I hate you so much." (Personal Attack, not a threat of violence)

Blackmail	Coercive	Must involve a	"Send	"You should
	demands	clear 'if you	nudes or	give me \$100."
	involving	don't X, then I	I'll tell	(demand, no
	threats to	will Y (reveal	everyone	threat)
	reveal	info/take	you	"I'm going to
	damaging,	action)'	cheated	tell the teacher
	sensitive, or	structure,	on the	you were late."
	private	explicit or	exam."	(reporting, not
	information	strongly	"If you	usually
	(true or	implied.	don't pay	blackmail
	false) or to		me \$100,	unless tied to a
	take other		I'm posting	coercive
	harmful		these	demand for
	action unless		photos of	something
	a specific		you."	unrelated)
	demand		you.	arri clatea)
	(e.g., for			
	money,			
	actions,			
	nudes) is			
	met.			

hostile, and	hominem	complete	dumb."
abusive	attacks,	idiot."	(criticism of
language,	name-calling,	"What a	idea, not
insults, or	and degrading		person)
slurs	language	pathetic loser you	"I disagree with
targeting an	directed at the	are."	you."
individual's	person. More	are.	(disagreement)
character,	severe than	"You're	(disagreement)
intelligence,	General Insult.	worthless."	"You're not
appearance		"You're	making sense."
(where not		such an	(can be
Body-Shamin		asshole."	borderline; if
g), or			hostile like
attributes.			"Nobody cares
This is about			about your
attacking the			opinion." (can
*person*,			be Exclusion or
not their			General Insult,
argument or			typically less
identity			aggressive than
(unless it			PA)
also qualifies			•
for			
Identity-Base			
d			

	Harassment)			
Identity-Based Harassment	Harassment, slurs, or attacks targeting someone based on their actual or perceived membership in a protected group (e.g., race, ethnicity, gender, religion, sexual orientation, disability, nationality).	The attack must explicitly or clearly implicitly reference a protected characteristic as the reason for the hostility.	"Go back to your country, [slur for nationality ]."  "Women are too emotional to lead." (directed at a woman in context)  "Typical [religious slur] behavior."	"You suck."  (Personal Attack, no identity link)  "Your music taste is terrible."  (opinion, not identity-based)  "He's acting crazy." (unless 'crazy' is used as a slur against someone with a known mental health condition in a targeted way)

Body-Shaming	Criticism,	Specific to	"You're so	"I hate your
	mockery, or	physical	fat you	outfit today."
	insults based	attributes and	need your	(criticism of
	on an	aims to	own zip	choice, not
	individual's	demean based	code."	body)
	physical	on them.	"Look at	"You're ugly."
	body shape,		those twig	(can be
	size, specific		arms,	Personal
	body parts,		pathetic."	Attack; if it
	or physical		patrictic.	clearly refers to
	appearance		"She has a	overall physical
	attributes		terrible	appearance in
	(e.g., "too		skin."	a shaming way,
	fat," "too			it's
	skinny," "ugly			Body-Shaming.
	nose").			If it's more
				general insult,
				PA). Prioritize
				Body-Shaming
				if specific
				physical traits
				are targeted for
				mockery.
				,

Exclusion	Statements or actions clearly intended to ostracize, isolate, socially reject, or make an individual feel unwelcome or unwanted in a group or conversation .	Focus on the act of pushing someone away or making them an outcast.	"Nobody likes you, just leave."  "You're not invited to this chat anymore."  "We'd all be better off if you weren't here."	"I don't agree with you." (disagreement)  "I need some space." (personal boundary, not necessarily exclusion from a group)
-----------	--	---	--	--

Sexual Harassment	Unwanted sexual comments, advances, objectificatio n, sexual propositions, or requests for sexual favors. This includes derogatory gendered sexual terms.	content must be clearly sexual in nature and unwanted (context may indicate unwelcomenes s).	"Show me your tits."  "Bet you'd be great in bed, [sexualize d term]."  "Are you a virgin?  Want me to change that?"	"You look nice today." (compliment, typically not SH unless context of repetition/pow er makes it so) "That dress is very flattering." (compliment) "Let's go on a date." (social invitation, not inherently SH unless coercive or repeatedly unwanted)
-------------------	---	--	--	---

Encouragement of Self-Harm	Directly urging, encouraging, or inciting someone to harm themselves physically, commit suicide, or engage in self-destructi ve behaviors (e.g., eating disorders).	Must be a clear call or suggestion for the target to inflict harm upon themselves.	"Go kill yourself."  "You should just cut yourself, you deserve it."  "The world would be better if you starved yourself."	"Are you feeling suicidal?" (inquiry, not encouragemen t)  "I hate you so much I wish you were dead." (expression of hate like General Insult, not direct encouragemen t for self-harm;)  "Life is pointless sometimes."
			starved	encouragemen t for self-harm;) "Life is pointless

General Insult Content that Rude, Apply ONLY if "Nobody (Fallback no other harm clearly fits into dismissive, cares Category) label fits. It or mildly about your another should pointless defined insulting story." language generally NOT category (e.g., that is clearly co-occur with "You're an "You're idiot" is a meant to be more specific being offensive or harm types. If a Personal really demeaning message is a Attack). annoying." but does not "Personal Content that "That's a meet the Attack," label it isn't harmful. stupid criteria for as such, not as thing to "Your any other GI. This is for say." argument is specific insults that are (borderline weak." harm less direct or , if said (criticism of less severe category mildly and argument, not above (e.g., than a not insult) not a direct "Personal attacking Attack." Personal core Attack on intelligenc character, e, could be not GI: if identity-base hostile, d, not a could be threat, etc.). PA) Use this as a

last resort	"I hate
for harmful	you, i wish
content.	you were
	dead"
	(while this
	is wishing
	death
	upon
	someone,
	it seems to
	be
	directing
	hostile
	language
	but not
	threat so
	GI is
	option
	best fitting
	the text)

No Harm	No	The default if	"Thank	— (any
	aggression,	no other label	you for	message that
	insult, or	clearly applies	your	fits a harm
	harmful	and the content	input."	category)
	content	isn't offensive.	III na la avia a	
	present.		"I'm having	
	Polite,		a great day!"	
	neutral,		uay:	
	positive, or		"This is an	
	benign		interesting	
	statements.		article."	
	Also use if			
	harm is			
	extremely			
	ambiguous			
	or relies			
	heavily on			
	assumptions			
	not			
	supported			
	by the			
	message or			
	context.			

SKIP	Use this	Your well-being	_	_
	option for	is a priority.		
	content that	Use when		
	is personally	needed. No		
	distressing	judgment.		
	or potentially			
	harmful to			
	your			
	well-being.			
	You are not			
	expected to			
	annotate it.			

Important Note on "General Insult":

"General Insult" is a fallback category. Only apply it if the message contains rude, demeaning, or insulting content that does not clearly meet the threshold or definition of any other specific harm category (like Personal Attack, Identity-Based Harassment, etc.). If a more specific harm category applies, use that instead of, or as a priority over, General Insult. For example, a direct attack like "You're an idiot" is a Personal Attack, not a General Insult. "Your opinion is worthless" might be a General Insult if it doesn't feel like a direct character attack but is still demeaning.

# Handling Overlapping Categories:

Some messages may legitimately contain elements that fit multiple categories. In such cases, select all specific harm categories that apply. Avoid using "General Insult" if all parts of the harmful message are covered by more specific labels.

Message	Applicable Labels	Reason
"You fat loser, go away before I make you."	Body-Shaming + Personal Attack + Exclusion + Threats of Violence	"fat" targets appearance (Body-Shaming), "loser" is a direct insult to character (Personal Attack), "go away" attempts to ostracize (Exclusion), and "before I make you" implies physical force (Threats of Violence).
"If you don't send me that \$50, I swear I'm going to hurt you and post those embarrassing photos of you everywhere."	Blackmail + Threats of Violence	"If you don't send me that \$50" coupled with threats to "hurt you" (Threats of Violence) and "post those embarrassing photos" (revealing damaging information) constitutes Blackmail. The physical threat is also explicitly a Threat of Violence.

"Get your filthy	Identity-Based	Contains a racial slur and	
[racial slur] hands off	Harassment +	generalization ("You people are	
that! You people are	Personal Attack +	all thieves") targeting identity	
all thieves. Go kill	Encouragement of	(Identity-Based Harassment).	
yourself."	Self-Harm	"Filthy" and "thieves" are also	
		direct attacks (Personal Attack).	
		"Go kill yourself" is direct	
		Encouragement of Self-Harm.	
"You're worthless	Personal Attack +	"Worthless" is a degrading	
and everyone knows	Encouragement of	statement directly attacking the	
it. You should just	Self-Harm	person (Personal Attack), and	
die."	Schridin	"You should just die"	
uie.		_	
		(Encouragement of Self-Harm).	
		encourages suicide (Encouragement of Self-Harm).	

## Decision-Making Flow:

- 1. Assess for Severe & Specific Harms: First, check for unambiguous instances of:
  - Threats of Violence
  - Blackmail
  - Encouragement of Self-Harm
  - Sexual Harassment
- 2. Assess for Identity-Based & Direct Personal Attacks:If the above are not present, or in addition to them, check for:

- Identity-Based Harassment
- Personal Attack (direct, abusive insults to character/person)
- Body-Shaming
- Exclusion
- 3. Consider General Insult (Fallback): If the message is clearly insulting, rude, or demeaning, but \*does not meet the criteria for any of the more specific categories above\*, then consider General Insult.
- 4. No Harm: If no harmful intent or content is reasonably perceived according to the definitions, select No Harm.
- 5. Multi-Label: Remember to apply all relevant specific labels if different parts of the message qualify for different harm types.

#### Handling Unclear or Borderline Messages:

- Ask: "Is the message \*explicitly\* or \*very clearly and strongly implying\* harm as defined in one of the categories?"
- Focus on what is directly stated or unmistakably implied by the language used.
   Avoid making assumptions about user history or relationships not evident from the provided text.
- If harm is present but genuinely ambiguous between two specific categories, and both seem plausible, you may apply both. However, strive for the most precise categorization.
- If the message is harsh but its classification as a specific harm type is unclear, or it's vaguely offensive without fitting a clear definition:
  - First, see if it fits General Insult (as a fallback).

- If it doesn't even meet the criteria for General Insult, or if the harmful nature itself is truly ambiguous (i.e., could a reasonable person interpret this as non-harmful in context?), lean towards No Harm.
- When in doubt between a specific harm and No Harm, and the implication of harm is not strong, prefer No Harm. Do not assume harm where it isn't obvious or strongly implied.

#### General Annotation Instructions:

- 1. Read Carefully: Review the Target Message and its preceding context thoroughly before assigning any label(s).
- Label Based on Definitions: Adhere strictly to the provided definitions and examples.
   Label based on the reasonably perceived meaning of the message, not your personal feelings or offense level.
- 3. Apply All Relevant Labels: If a message contains multiple distinct types of harm, apply all applicable labels.
- 4. Avoid Assumptions: Only label content based on explicit statements or strong, direct implications within the provided text. Do not infer harm based on possibilities or what might be true outside the conversation.
- Consistency is Key: Refer back to this guide, especially the definitions, differentiators, and examples, to ensure consistent labeling across similar messages.
- 6. Sensitive Content & SKIP: If you encounter highly disturbing material that affects your well-being, use the SKIP option.
- 7. Take Breaks: Annotation, especially of sensitive content, can be demanding. Take regular breaks to maintain focus, accuracy, and well-being.