## **Appropriate Filtering for Education settings**



### **May 2025**

### **Filtering Provider Checklist Reponses**

Schools (and registered childcare providers) in England and Wales are required "to ensure children are safe from terrorist and extremist material when accessing the internet in school, including by establishing appropriate levels of filtering". Furthermore, it expects that they "assess the risk of [their] children being drawn into terrorism, including support for extremist ideas that are part of terrorist ideology". There are a number of self review systems (eg www.360safe.org.uk) that will support a school in assessing their wider online safety policy and practice.

The Department for Education's statutory guidance 'Keeping Children Safe in Education' obliges schools and colleges in England to "ensure appropriate filters and appropriate monitoring systems are in place and regularly review their effectiveness" and they "should be doing all that they reasonably can to limit children's exposure to [Content, Contact, Conduct, Contract] risks from the school's or college's IT system" however, schools will need to "be careful that "over blocking" does not lead to unreasonable restrictions as to what children can be taught with regards to online teaching and safeguarding."

By completing all fields and returning to UK Safer Internet Centre (<a href="mailto:enquiries@saferinternet.org.uk">enquiries@saferinternet.org.uk</a>), the aim of this document is to help filtering providers to illustrate to education settings (including Early years, schools and FE) how their particular technology system(s) meets the national defined 'appropriate filtering standards. Fully completed forms will be hosted on the UK Safer Internet Centre website alongside the definitions

It is important to recognise that no filtering systems can be 100% effective and need to be supported with good teaching and learning practice and effective supervision.

Company / Organisation	Talk Straight / Schools Broadband		
Address	Units 2 – 4 Backstone Business Park, Dansk Way, Ilkley, LS29 8JZ		
Contact details	info@talk-straight.com		
Filtering System	Netsweeper / Incident Management Platform		
Date of assessment	20/10/2025		

#### System Rating response

Where a supplier is able to confirm that their service fully meets the issue identified in a specific checklist the appropriate self-certification colour for that question is GREEN.	
Where a supplier is not able to confirm that their service fully meets the issue	
identified in a specific checklist question the appropriate self-certification colour	
for that question is AMBER.	

.

## **Illegal Online Content**

Filtering providers should ensure that access to illegal content is blocked, specifically that the filtering providers:

Aspect	Rating	Explanation
<ul> <li>Are IWF members</li> </ul>		Yes, we have been members
		of the IWF for over 14 years
<ul> <li>and block access to illegal Child Abuse</li> </ul>		Yes, our Netsweeper service
Images (by actively implementing the IWF		integrate with the IWF CAIC
URL list), including frequency of URL list		illegal content list. The IWF
update		functionality is not exposed
		in the graphical user
		interface and cannot be
		disabled.
<ul> <li>Integrate the 'the police assessed list of</li> </ul>		Yes, our Netsweeper service
unlawful terrorist content, produced on		integrate this list into their
behalf of the Home Office'		Web Filtering service and
		this is blocked by default.
<ul> <li>Confirm that filters for illegal content cannot</li> </ul>		Yes, all categories associated
be disabled by anyone at the school		with illegal content are
(including any system administrator).		locked at the system level
		and all school administrators
		cannot enable these
		categories.

Describing how, their system manages the following illegal content

Content	Explanatory notes – Content that:	Rating	Explanation
child sexual	Content that depicts or promotes		Netsweeper has its own CSAM
abuse	sexual abuse or exploitation of		(Child sexual Abuse Material)
	children, which is strictly		category within the platform
	prohibited and subject to severe		which it uses to actively block this
	legal penalties.		harmful content. On top of this
			Netsweeper pulls URLS directly
			from IWF and Project Arachnid.
controlling or	Online actions that involve		Netsweeper covers this area with
coercive	psychological abuse, manipulation,		multiple categories to cover
behaviour	or intimidation to control another		content related to Controlling
	individual, often occurring in		and coercive behaviour.
	domestic contexts.		
extreme	Content that graphically depicts		Netsweeper categorisation allows
sexual	acts of severe sexual violence,		multiple categorises to be used
violence	intended to shock or incite similar		against URLs meaning that
	behaviour, and is illegal under UK		Extreme Sexual Violence would
	law.		be cover by the following
			Categorises – Violence, Extreme
			and Adult content. This granular
			nature ensures URLs do not get
			miscategorised due to ridged
			categorisation. It also ensures
			that DSLs are aware of the type

		of content that has been
		accessed when reporting.
extreme	Pornographic material portraying	Netsweeper categorisation allows
pornography	acts that threaten a person's life or	multiple categorises to be used
	could result in serious injury, and is	against URLs meaning that
	deemed obscene and unlawful.	Extreme Pornography would be
		cover by the following
		Categorises – Extreme and
		Pornography. This granular
		nature ensure URLs do not get
		miscategorised due to ridged
		categorisation. It also ensures
		that DSLs are aware of the type
		of content that has been
		accessed when reporting.
fraud	Deceptive practices conducted	Netsweeper's Scam Category
	online with the intent to secure	cover all scam and fraud related
	unfair or unlawful financial gain,	content. Netsweeper also works
	including phishing and scam	with organisations like Scam
	activities.	advisors to block known
		Fraudulent activity.
racially or	Content that incites hatred or	Netsweeper has multiple
religiously	violence against individuals based	categories to cover this subject
aggravated	on race or religion, undermining	matter including Profanity and
public order	public safety and cohesion.	hate speech. Netsweeper would
offences	,	label the URL based on its specific
		type of offence.
inciting	Online material that encourages or	Netsweeper categorisation allows
violence	glorifies acts of violence, posing	multiple categorises to be used
	significant risks to public safety	against URLs meaning that
	and order.	inciting violence would be cover
		by the following Categorises –
		Extreme, Hate speech and
		violence. This granular nature
		ensure URLs do not get
		miscategorised due to ridged
		categorisation. It also ensures
		that DSLs are aware of the type
		of content that has been
		accessed when reporting.
illegal	Content that promotes or	Netsweeper Criminal skills
immigration	facilitates unauthorized entry into	category would capture all
and people	a country, including services	content related to both illegal
smuggling	offering illegal transportation or	immigration and people
<b>-</b>	documentation.	smuggling
promoting or	Material that encourages or assists	Netsweeper Self harm and
facilitating	individuals in committing suicide,	suicide category covers all
suicide	posing serious risks to vulnerable	content related to this matter.
	populations.	Netsweeper also works with
		leading Suicide Prevention charity
		"R;pple" to keep on top of all

		new risk and word used by users looking at suicide ideation content.
intimate image abuse	The non-consensual sharing of private sexual images or videos, commonly known as "revenge porn," intended to cause distress or harm.	Netsweeper can block sexual content via its Adult and pornography categories. On top of this, methods of distribution like personal mail box, communication tools (outside schools control) or storage tools are covered by Web email, Web storage and web chat.
selling illegal drugs or weapons	Online activities involving the advertisement or sale of prohibited substances or firearms, contravening legal regulations.	Netsweeper's Substance abuse, marijuana, weapons, criminal skills and sale would categorise this content. As an example, an online storing selling drugs would be categorised as – Substance abuse, Criminal Skills and Sales due to the nature of the content. This level of detail ensure DSL are aware of the content accessed by their end users.
sexual exploitation	Content that involves taking advantage of individuals sexually for personal gain or profit, including trafficking and forced prostitution.	Netsweeper following categories would cover this – Adult, Pornography and criminal skills.
Terrorism	Material that promotes, incites, or instructs on terrorist activities, aiming to radicalise individuals or coordinate acts of terror.	Netsweeper has a dedicated Terrorism category for all terrorism related content.

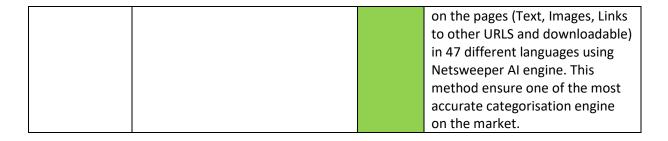
## **Inappropriate Online Content**

Recognising that no filter can guarantee to be 100% effective, providers should both confirm, and describe how, their system manages the following content

Content	Explanatory notes – Content that:	Rating	Explanation
Gambling	Enables gambling		This includes sites that encourage or provide information on the wagering or risking of money or any valuables on a game, contest, or other event in which the outcome is partially or completely dependent upon chance or on one's abilities. Sites that promote or facilitate gambling information, as well those that are purely factual and strategic sites that promote

Hate speech / Discriminiatio n	Content that expresses hate or encourages violence towards a person or group based on something such as disability, race, religion, sex, or sexual orientation. Promotes the unjust or prejudicial treatment of people with protected characteristics of the Equality Act 2010	cheating are also included. The excludes sites that are clearly support sites for gambling addiction as well as travel destination sites that do not enable gambling  Netsweeper Hate Speech category identifies content related to hate speech and discrimination content. Using Netsweeper's real time categorisation engine (used by Netsweeper for over 25 years) Netsweeper analysis the conte on the pages (Text, Images, Lit to other URLS and downloada in 47 different languages using Netsweeper AI engine. This method ensures one of the me accurate categorisation engine on the market.	y ) ent nks ble) g
Harmful content	Content that is bullying, abusive or hateful. Content which depicts or encourages serious violence or injury. Content which encourages dangerous stunts and challenges; including the ingestion, inhalation or exposure to harmful substances.	Netsweeper's Bullying, Violence and extreme category identify content related to harmful content. Using Netsweeper's retime categorisation engine (use by Netsweeper for over 25 years) Netsweeper analysis the content on the pages (Text, Images, Linto other URLS and downloads in 47 different languages using Netsweeper AI engine. This method ensures one of the metaccurate categorisation engine on the market.	real sed ars) ent nks ble)
Malware / Hacking	promotes the compromising of systems including anonymous browsing and other filter bypass tools as well as sites hosting malicious content	Netsweeper has an extensive malicious content database us by our customers. Covered by Netsweeper categories Malwa Phishing, Virus, Infected Host, and on top this a range of som of the top leading security vendors list that we managed and maintain to safeguard our customers networks.	are, ne
Mis / Dis Information	Promotes or spreads false or misleading information intended to deceive, manipulate, or harm, including content undermining trust in factual information or institutions	Netsweeper News category ca be used to block website that known to promote mis and dis information. This walled garde approach means that trusted	are s

Piracy and copyright theft	includes illegal provision of copyrighted material	new authorities can be allowed on the school network.  Netsweeper Copyright Infringement Category blocks content contains sites that use, provide, or distribute information on copyrighted intellectual property or illicitly copied material which violates the owners' rights. On top of this Netsweeper uses the PIPCU (Police intellectual Property Crime Unit) list which guards against Piracy and copyright offences.
Pornography	displays sexual acts or explicit images	Netsweeper's Pornography category identifies content related to harmful content. Using Netsweeper's real time categorisation engine (used by Netsweeper for over 25 years) Netsweeper analysis the content on the pages (Text, Images, Links to other URLS and downloadable) in 47 different languages using Netsweeper Al engine. This method ensure one of the most accurate categorisation engine on the market.
Self Harm and eating disorders	content that encourages, promotes, or provides instructions for self harm, eating disorders or suicide	Netsweeper's Self harm category identifies content related to harmful content. Using Netsweeper's real time categorisation engine (used by Netsweeper for over 25 years) Netsweeper analysis the content on the pages (Text, Images, Links to other URLS and downloadable) in 47 different languages using Netsweeper Al engine. This method ensures one of the most accurate categorisation engine on the market.
Violence Against Women and Girls (VAWG)	Promotes or glorifies violence, abuse, coercion, or harmful stereotypes targeting women and girls, including content that normalises gender-based violence or perpetuates misogyny.	Netsweeper's Hate Speech, Violence and Criminal Skills categories identify content related to harmful content. Using Netsweeper's real time categorisation engine (used by Netsweeper for over 25 years) Netsweeper analysis the content



This list should not be considered an exhaustive list. Please outline how the system manages this content and many other aspects

The Netsweeper collective community numbers are over billion devices worldwide. This collective, together with our AI technology and human oversight defines URL classification. Netsweeper publishes classification of filtering and categorises on the Netsweeper website as well as a view in real time of new content categorised. Netsweeper's core competency is using our patented techniques to categorise every URL that passes through our deployed systems. Netsweeper is both real-time and employs a hierarchy of data (URL-to-category), with our Category Naming Service (CNS) as a global-master database. If any customer anywhere in the world accesses a URL, that URL is submitted to the local policy server, if that policy server cannot find a category match, it is automatically submitted to the CNS and looked up there. If the CNS already has the category mapping it is immediately returned to the local system and cached there for future use, a policy decision is then made by the policy server. If neither the local system, nor the CNS has a category match, the URL is submitted to our "Artificial Intelligence" system that will interrogate the content at-and-around that URL, assess the content, detect if it references or contains malware, and assigns one or more categories to the URL into the CNS and then back to the local system, a policy decision is then made by the policy server. The CNS allows us to adapt to trending URLs immediately due to its world-wide scope. If the local system hasn't seen a particular URL yet, then CNS probably has. If the URL has been assigned one or more categories, local systems see immediate responses (sub-second). If the URL is truly "new" then the AI will typically process the content within 4-5 seconds. The local policy servers can be configured with techniques to minimise the "new URL" wait.

Regarding the duration and extent of logfile (Internet history) data retention, providers should outline their retention policy, specifically including the extent to the identification of individuals and the duration to which all data is retained.

We typically retain logs on our deployment for 24 months, however schools, should they wish can choose to have longer retention.

Providers should be clear how their system does not over block access so it does not lead to unreasonable restrictions

Netsweeper policy manager offers granular controls around categorisation. School leaders can enforce the policies based on the guidance set out by the DFE and the Uk safer internet centre. Netsweeper Realtime categorisation ensures the ever-changing landscape of the internet is kept on top of by Netsweeper reviewing content on a regular basis. This ensures categories change when the content changes. Netsweeper's ability to label URLs with multiple categories ensures Netsweeper can be more specific on categorisation ensuring schools don't need to over block content due to the risk associated with non-granular categorisation.

## Filtering System Features

How does the filtering system meet the following principles:

Principle	Rating	Explanation
<ul> <li>Context appropriate differentiated filtering, based on age, vulnerability and risk of harm – also includes the ability to vary filtering strength appropriate for staff</li> </ul>		Netsweeper allows for differentiated filtering based on different user types. For example, students may receive differentiated filtering based on age.  Vulnerable users may have certain categories blocked that are not blocked for other uses
Circumvention – the extent and ability to identify and manage technologies and techniques used to circumvent the system, specifically VPN, proxy services, DNS over HTTPS and ECH.		Netsweeper uses advanced technology to detect and prevent attempts to circumvent the system. Talk Straight have our own DNS servers to specifically stop ECH requests. Our hosted FortiGate firewall services can also be setup to further inspect traffic and stop VPN and proxy services in conjunction with Netsweeper filtering and blocked domain lists.
Control – has the ability and ease of use that allows schools to control the filter themselves to permit or deny access to specific content. Any changes to the filter system are logged enabling an audit trail that ensure transparency and that individuals are not able to make unilateral changes		All changes made to the system is logged and auditable by the school. The schools are able to administer their own content filtering. For MATs, delegated access can also be created so MAT administrators could access all schools setups but local DSLs could only change / view data applicable to their school. Every request is logged for auditing and compliance purposes.
<ul> <li>Contextual Content Filters – in addition to URL or IP based filtering, Schools should understand the extent to which (http and https) content is dynamically analysed as it is</li> </ul>		Netsweeper provides on the fly categorisation based on the content and context of text and

streamed to the user and blocked. This would include AI or user generated content, for example, being able to contextually analyse text and dynamically filter the content produced (for example ChatGPT). For schools' strategy or policy that allows the use of AI or user generated content, understanding the technical limitations of the system, such as whether it supports real-time filtering, is important.	links that appear on the page. Netsweeper Al category can also be used to safeguard users against more harm generative AI tools.
Deployment – filtering systems can be deployed in a variety (and combination) of ways (eg on device, network level, cloud, DNS). Providers should describe how their systems are deployed alongside any required configurations	Our Netsweeper service is delivered via a Cloud only model. Network Level and/or client lead. Netsweeper filtering provides device level filtering both on network and off network across Windows, Chromebooks, Mac, iPad, iPhones and Android device. This ensures the most comprehensive filtering deployment.
Filtering Policy – the filtering provider publishes a rationale that details their approach to filtering with classification and categorisation as well as how the system addresses over blocking	Netsweeper broad categorisation engine ensures flexibility when applying policies ensure school traffic isn't over blocked.
<ul> <li>Group / Multi-site Management – the ability for deployment of central policy and central oversight or dashboard</li> </ul>	Netsweeper multi tenanted platform offers schools and trusts the ability to deliver granular policy management both at a per user level and per school basis.
<ul> <li>Identification - the filtering system should have the ability to identify users and devices to attribute access (particularly for mobile devices) and allow the application of appropriate configurations and restrictions for individual users. This would ensure safer and more personalised filtering experiences.</li> </ul>	Users can be identified by various methodologies including but not limited to Entra (formally Azure), Local AD, google directory. Netsweeper supports hybrid and cross directory services (School that have both Entra and Google tenancy).
Mobile and App content – mobile and app  content is often delivered in entirely different.	Netsweeper is capable of
content is often delivered in entirely different	filtering all device based

mechanisms from that delivered through a	content including
traditional web browser. To what extent does	content delivered via
the filter system block inappropriate content	apps.
via mobile and app technologies (beyond	
typical web browser delivered content).	
Providers should be clear about the capability	
of their filtering system to manage content on	
mobile and web apps and any configuration or	
component requirements to achieve this	
<ul> <li>Multiple language support – the ability for the</li> </ul>	Netsweeper performs
system to manage relevant languages	dynamic filtering in 47
	languages
Remote devices – with many children and	Netsweeper can deliver
staff working remotely, the ability for school	identical filtering to the
owned devices to receive the same or	school's devices onsite
equivalent filtering to that provided in school	and offsite at a device
	level (not just browser
	based) across Windows,
	Chromebooks, Mac,
	iPad, iPhone and
	Android device
<ul> <li>Reporting mechanism – the ability to report</li> </ul>	Our enhanced Incident
inappropriate content for access or blocking	Management Reporting
mappings contain in access of closiming	Platform mechanism
	allows DSL the ability to
	report on inappropriate
	content accessed or
	blocked both on and off
	network. Reports and
	alerts can be delivered
	via a number of
	mechanisms including
	email, teams and slack.
Reports – the system offers clear granular	Our Incident
historical information on the websites users	Management Platform
have accessed or attempted to access	can provide detailed
have accessed of attempted to access	reports on with the who,
	what, when and where.
	•
	Our templated reports
	and custom reports can
	be set up to proactively
	send known
	safeguarding threats to
	DSL in near real time.
<ul> <li>Safe Search – the ability to enforce 'safe</li> </ul>	Safe search can be
search' when using search engines	enforced across all well
	known search engines
	and lesser-known ones.
<ul> <li>Safeguarding case management integration –</li> </ul>	Netsweeper integrates
the ability to integrate with school	with case management
	solutions either via

direct API to the platform or CSV upload.

# How does your filtering system manage access to Generative AI technologies (e.g. ChatGPT, image generators, writing assistants)?

In your response, please describe whether and how your system identifies, categorises, or blocks Generative AI tools; how access can be controlled based on age, risk, or educational need; any limitations in filtering AI-generated content—particularly where such content is embedded within other platforms or applications; and what support or configuration guidance you offer to schools to help them align with the UK Safer Internet Centre's Appropriate Filtering Definitions and relevant national safeguarding frameworks.

Netsweeper's AI categorisation engine can help identify generative AI technology. Netsweeper would suggest using the walled garden approach to safeguarding users from generative AI. As an example schools should only use trust tools like Co-Pilot inside a school's tenancy, thus blocking all over technologies.

Filtering systems are only ever a tool in helping to safeguard children when online and schools have an obligation to "consider how children may be taught about safeguarding, including online, through teaching and learning opportunities, as part of providing a broad and balanced curriculum".<sup>1</sup>

Please note below opportunities to support schools (and other settings) in this regard

Ripple – Suicide prevention charity – www.ripplesuicideprevention.com

\_

<sup>&</sup>lt;sup>1</sup> https://www.gov.uk/government/publications/keeping-children-safe-in-education--2

#### **PROVIDER SELF-CERTIFICATION DECLARATION**

In order that schools can be confident regarding the accuracy of the self-certification statements, the supplier confirms:

- that their self-certification responses have been fully and accurately completed by a person or persons who are competent in the relevant fields
- that they will update their self-certification responses promptly when changes to the service or its terms and conditions would result in their existing compliance statement no longer being accurate or complete
- that they will provide any additional information or clarification sought as part of the selfcertification process
- that if at any time, the UK Safer Internet Centre is of the view that any element or elements of a provider's self-certification responses require independent verification, they will agree to that independent verification, supply all necessary clarification requested, meet the associated verification costs, or withdraw their self-certification submission.

Name	David Tindall
Position	CEO
Date	20/10/2025/
Signature	DAVID TINDALL